

Combining inverse photogrammetry and BIM for automated labeling of construction site images for machine learning

Alex Braun^{a,b}, André Borrmann^{a,b}

^a*Chair of Computational Modeling and Simulation, Technical University of Munich, Germany*

^b*Leonhard Obermeyer Center, TUM Center of Digital Methods for the Built Environment*

Abstract

Image-based object detection provides a valuable basis for site information retrieval and construction progress monitoring. Machine learning approaches, such as neural networks, are able to provide reliable detection rates. However, labeling of training data is a tedious and time-consuming process, as it must be performed manually for a substantial number of images. The paper presents a novel method for automatically labeling construction images based on the combination of 4D Building Information Models and an inverse photogrammetry approach. For the reconstruction of point clouds, which are often used for progress monitoring, a large number of pictures are taken from the site. By aligning the Building Information Model and the resulting point cloud, it is possible to project any building element of the BIM model into the acquired pictures. This allows for automated labeling as the semantic information of the element type is provided by the BIM model and can be associated with the respective regions. The labeled data can subsequently be used to train an image-based neural network. Since the exact regions for all elements are defined, labels can be generated for basic tasks like classification as well as more complex tasks like semantic segmentation. To prove the feasibility of the developed methods, the labeling procedure is applied to several real-world construction sites, providing over 30,000 automatically labeled elements. The correctness of the assigned labels has been validated by pixel based area comparison against manual labels.

Keywords: Machine Learning, Labeling, Construction progress monitoring, BIM, Photogrammetry, semantic and temporal knowledge

1. Introduction

Large construction projects require a variety of different manufacturing companies of several trades on site (for example masonry, concrete and metal works, HVAC, ...). An important goal for the main contractor is to keep track of accomplished tasks by subcontractors to maintain the general schedule. Additionally, the documentation of correctly executed tasks plays a crucial role for all involved parties. In construction, process supervision and monitoring is still a mostly analog and manual task. To prove that the work has been completed as defined per contract, all performed tasks have to be monitored and documented. The demand for a complete and detailed monitoring technique rises for large construction sites where the complete construction area becomes too large to monitor by hand, and the number of subcontractors rises. Main contractors that control their subcontractors' work need to keep an overview of the current construction state. Regulatory issues add up on the requirement to keep track of the current status on site.

The ongoing digitization and the establishment of building information modeling (BIM) technologies in the planning of construction projects help to establish new methods for process optimization. In an ideal implementation of the BIM concept, all semantic data on materials, construction methods, and even the process schedule are connected. On this basis, it is possible to make much more precise estimations about the project costs and its duration. Most importantly, possible deviations from the schedule can be detected early, and the resources can be adapted accordingly.

This technological advancement allows new methods in construction monitoring. In Braun et al. [1], the authors propose a method for automated progress monitoring using photogrammetric point clouds and 4D Building Information Models. The central concept is to use standard camera equipment on construction sites to capture the current construction state by taking pictures of the complete facility under construction at regular intervals. As soon as a sufficient number of images from different points of view are available, a 3D point cloud can be reconstructed with the help of photogrammetric methods. This point cloud represents one particular time-stamp of the construction progress (as-built) and is subsequently matched against the geometry of the BIM (as-planned) on a per-element basis.

Figure 1 shows the C#-based WPF software tool, developed in the scope of this research. The tool visualizes a building information model and all corresponding semantic data. Additionally, the observation results can be

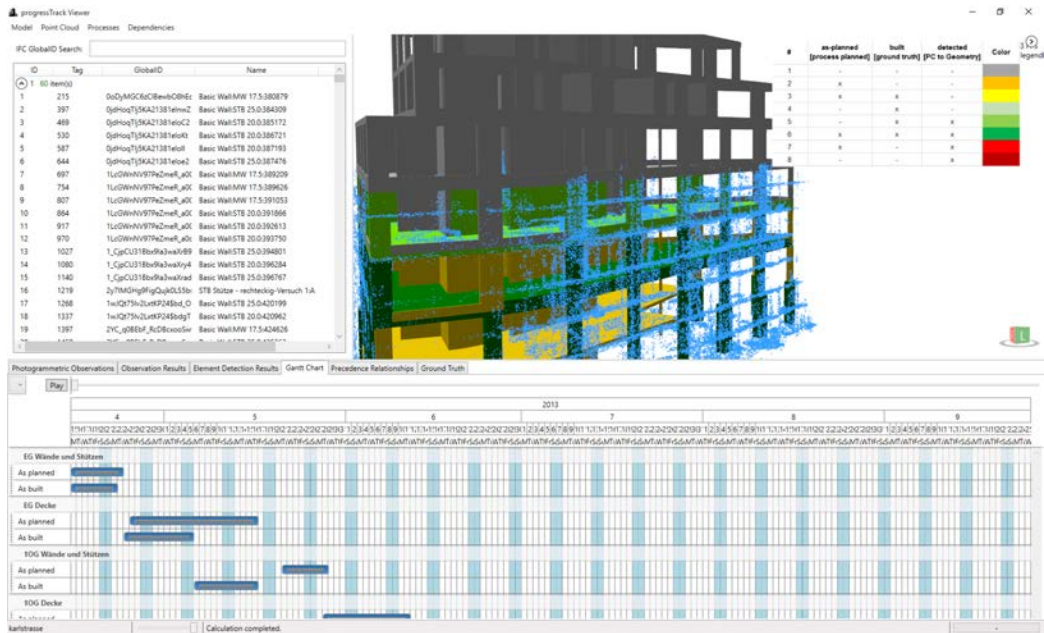


Figure 1: progressTrack: 4D BIM viewer incorporating detection states, process information and point clouds from observations

38 selected and are supported by the possible overlay of the corresponding point
 39 clouds.

40 The presented approach can be varied in terms of acquisition method
 41 (laser scanning, manual acquisition, ...) and matching methods (as dis-
 42 cussed in Section 2 - Related work). However, none of the methods is capable
 43 of providing absolute reliability due to occlusions or other boundary condi-
 44 tions. To further improve the reliability of the methods mentioned above,
 45 image-based machine learning techniques offer a promising approach. These
 46 techniques allow to analyze pictures based on their contents and even mark
 47 and classify specific regions of pictures. This new information can further
 48 improve the geometric as-planned vs. as-built comparison based on point
 49 clouds by increasing the reliability of made assumptions while comparing
 50 semantic data from the BIM with classified categories on similar pictures.

51 Recently, Convolutional Neural Networks (CNN) were introduced in this
 52 context [2, 3]. These networks require large training sets to learn similari-
 53 ties of provided data-sets to make assumptions on unknown data. Applica-
 54 tions of CNNs range from face-detection in security-related applications to

55 autonomous driving [4]. With respect to automated construction monitor-
56 ing, these methods can help to detect construction elements on pictures and
57 provide an alternative method for detection in case of low point cloud den-
58 sities and to improve the overall accuracy of detection [5, 6]. However, data
59 pre-processing and labeling of test-sets for the training of said algorithms
60 is a laborious and time-consuming task since common CNNs require large
61 amounts of labeled data [7].

62 This paper presents a method to automate the process of construction-site
63 image labeling. The proposed method makes use of available information on
64 image localization from the photogrammetric process as well as information
65 on the presence of individual construction elements from the as-planned vs.
66 as-built comparison by the process described above. The resulting availability
67 of training data provides the basis for applying the trained CNN for image-
68 based object detection on any construction site, in particular, those where
69 a 4D-BIM does not exist or only a limited number of images are taken, and
70 the generation of a point cloud is not possible. However, this paper does not
71 report on these next stages but focuses on the provision of correctly labeled
72 images as an essential first step.

73 **2. Related work**

74 *2.1. Automated construction monitoring*

75 Several methods for BIM-based progress monitoring have been developed
76 in recent years [8]. Basic methods make use of minor technical advancements
77 like introducing email and tablet computers into the manual monitoring pro-
78 cess. These methods still require manual work, but already contribute to the
79 shift towards digitization. More advanced methods try to track individual
80 building components through radio-frequency identification (RFID) tags or
81 similar methods (for example QR codes).

82 Current state-of-the-art procedures apply vision-based methods for more
83 reliable element identification. These methods either make direct use of pho-
84 tographs or videos taken on site as input for image recognition techniques
85 or apply laser scanners or photogrammetric methods to create point clouds
86 that hold point-based 3D information and additionally color information.

87 Bosche and Haas [9], Bosché [10] present a system for as-planned vs. as-
88 built comparisons based on laser-scanning data. The generated point clouds
89 are co-registered with the model using an adapted Iterative-Closest-Point-
90 Algorithm (ICP). Within this system, the as-planned model is converted

91 into a point cloud by simulating the points using the known positions of
92 the laser scanner. For verification, they use the percentage of simulated
93 points, which can be verified by the real laser scan. Turkan et al. [11] use
94 and extend this system for progress tracking using schedule information for
95 estimating the progress in terms of earned value and for detecting secondary
96 objects. Kim et al. [12] detect specific component types using a supervised
97 classification based on Lalonde features derived from the as-built point cloud.
98 An object is regarded as detected if the type matches the type present in the
99 model. As above, this method requires that the model is sampled into a point
100 representation. Zhang and Arditi [13] introduce a measure for deciding four
101 cases (object not in place, point cloud represents a full object or a partially
102 completed object or a different object) based on the relationship of points
103 within the boundaries of the object and the boundaries of the shrunk objects.
104 The authors test their approach in a very simplified artificial environment,
105 which is significantly less challenging than the processing of data acquired
106 on real construction sites.

107 In comparison with laser scanning, photogrammetric methods are less
108 accurate. However, standard cameras have the advantage that they can be
109 used more flexibly, and their costs are much lower. This leads to the need
110 for other processing strategies when image data is used. Omar and Nehdi [8]
111 give an overview and comparison of image-based approaches for monitoring
112 construction progress. Ibrahim et al. [14] use a single camera approach and
113 compare images taken during a specified period and rasterize them. The
114 change between two time-frames is detected using a spatial-temporal deriva-
115 tive filter. This approach is not directly bound to the geometry of a BIM and
116 therefore cannot identify additional construction elements on site. Kim et al.
117 [15] use a fixed camera and image processing techniques for the detection
118 of new construction elements and the update of the construction schedule.
119 Since many fixed cameras would be necessary to cover a whole construction
120 site, more approaches rely on images from hand-held cameras covering the
121 whole construction site.

122 For finding the correct scale of the point cloud, stereo-camera systems can
123 be used, as done in [16, 17, 18]. Rashidi et al. [19] propose using a colored
124 cube of known size as a target, which can be automatically measured to
125 determine the scale. Additionally, image-based approaches can be compared
126 with laser-scanning results [20]. The artificial test data is strongly simplified,
127 and the real data experiments are limited to a small part of a construction
128 site. Only relative accuracy measures are given since no scale was introduced

129 to the photogrammetry measurements. Golparvar-Fard et al. [21, 22] use
130 unstructured images of a construction site to create a point cloud. The
131 orientation of the images is computed using a Structure-from-Motion process
132 (SFM). Subsequently, dense point clouds are calculated. For the comparison
133 of as-planned and as-built geometry, the scene is discretized into a voxel
134 grid. The construction progress is determined in a probabilistic approach, in
135 which the threshold parameters for detection are determined by supervised
136 learning. This framework makes it possible to take occlusions into account.
137 This approach relies on the discretization of space as a voxel grid to the
138 size of a few centimeters. In contrast, the approach presented here is based
139 on calculating the deviation between a point cloud and the building model
140 directly and introduces a scoring function for the verification process.

141 The mentioned approaches provide valuable enhancements for automated
142 construction progress monitoring. However, so far, not all potential benefits
143 from using semantic BIM data are unlocked to their full extent. Also, current
144 research does not present solutions for occluded elements as well as temporary
145 construction elements like scaffolds. These elements cover large parts of
146 construction sites and thus cannot be neglected. The presented approach
147 tries to solve this issue by analyzing the images taken during the SFM process.

148 *2.2. Computer Vision*

149 Computer Vision is a heavily researched topic, that got even more atten-
150 tion through recent advances in autonomous driving and machine learning
151 related topics. Image analysis for construction sites, on the other hand, is
152 a rather new topic. Since one of the key aspects of machine learning is the
153 collection of large data-sets, current approaches focus on data gathering. In
154 the scope of automated progress monitoring, Han and Golparvar-Fard [23]
155 published an approach for labeling based on the commercial service Amazon
156 Turk. Chi and Caldas [24] used first versions of neural networks to detect
157 construction machinery on images, Kropp et al. [25] tried to detect in-door
158 construction elements based on similarities, focusing on radiators. Kim et al.
159 [26] used ML-based techniques for construction progress monitoring. They
160 analyzed images by filtering them to remove noise and uninteresting ele-
161 ments to focus the comparison on relevant construction processes. Other
162 publications mainly focus on defect detection (like for example cracks) in
163 construction images [27].

164 Current research mainly uses manual labels for computer vision. Addi-
165 tionally, no construction data set is currently covering the whole amount of

166 construction elements. An automated labeling approach could better this
167 lack of data to further improve machine learning methods in this scope of
168 application.

169 **3. Problem statement**

170 Monitoring of construction sites by applying photogrammetric methods
171 has become a common practice. Currently, several companies (for example
172 Pix4D, DroneDeploy) provide commercial solutions for end users that allows
173 to generate 3D meshes and point clouds from UAV-based site observations.
174 All these methods give reasonable solutions for finished construction sites or
175 visible elements of interest.

176 However, there are still many unsolved problems in monitoring construc-
177 tion sites. Photogrammetric methods are sensitive to low structured surfaces
178 or windows. Because of the used method, each element needs to be visible
179 from multiple (at least two) different points of view. Thus, elements inside of
180 a building cannot be reconstructed as they are not visible from a UAV flying
181 outside of the building. Monitoring inside a building is currently still under
182 heavy research [28] and not available in an automated manner as orientation
183 and observation in such mutable areas like construction sites is hard to tackle.
184 These problems lead to holes or misaligned points in the final point cloud,
185 that hinder accurate and precise detection of building elements. On the other
186 hand, laser scanning requires many acquisition points and takes significantly
187 more time and manual effort for acquisition. Finally, both techniques remain
188 with occlusions for regions that are not visible during construction.

189 As can be seen in Figure 2, another problem is elements that are occluded
190 by temporary construction elements. Especially scaffolds and formwork ele-
191 ments occlude the view on walls or slabs, making it harder for algorithms to
192 detect the current state of construction progress.

193 This paper proposes a method that is meant to overcome some of the lim-
194 itations of the available methods. It contributes to the final goal of exploiting
195 images as an information source for construction state detection, either as ad-
196 ditional information in case one of the methods mentioned above is applied,
197 or even as sole and primary information if a 4D BIM does not exist or an
198 insufficient number of images is available for photogrammetric detection. To
199 achieve this, the authors propose to apply CNNs for automated object detec-
200 tion. However, a huge set of correctly labeled images is required for training



Figure 2: Occluded construction elements in generated point cloud caused by scaffolding, formworks, existing elements and missing information during the reconstruction process

201 the CNN and achieve high precision and low recall. So far, the labeling pro-
202 cess had to be performed manually in a laborious and error-prone process.
203 This is why the authors propose to automate this process by making use of
204 the methods they originally developed for construction progress monitoring.
205 In particular, we use image localization from the photogrammetric process
206 as well as information on the presence of individual construction elements
207 from the as-planned vs. as-built comparison. This results in the availability
208 of the required high quality, high volume training data.

209 4. Automated labeling of images

210 An essential part of progress monitoring is the detection of an element's
211 status, i.e. to decide whether an element is still under construction (e.g.,
212 surrounded by formwork) or finished. Pure point-cloud-to-model matching
213 methods are facing difficulties in this regard as temporary and auxiliary con-
214 structions (such as formwork) usually are not included in the BIM model. As
215 proposed in Braun et al. [29], computer vision based methods can help here
216 and significantly improve the reliability of as-planned vs. as-built compari-
217 son. The basic idea is to use visual information to decide upon an element's
218 visibility status.

219 The authors propose the use of machine learning (ML) methods for image-
220 based detection of a construction element’s status. However, ML techniques
221 require a large set of labeled images for training. As currently large labeled
222 sets of construction site images or not available, the labeling has to be per-
223 formed manually in a tedious and time-consuming process. Generating these
224 labels automatically can drastically reduce preparation efforts for training
225 and improving such networks.

226 The proposed concept of automatic labeling is based on fusing information
227 available from the photogrammetric process (images and relative position of
228 the camera) and the information available from the 4D BIM (object type,
229 object position). Since the BIM and the resulting point cloud are aligned,
230 each BIM element can be projected onto the image initially taken for the
231 photogrammetric process. This allows to precisely identify the region covered
232 by a building element on a picture.

233 However, there is a significant problem remaining: Information on the ac-
234 tual presence of the element cannot be reliably taken from the 4D as-planned
235 BIM, as execution time very often deviates from the original schedule (which
236 is the underlying rationale for applying progress monitoring). At this point,
237 we benefit from the original point-cloud vs. BIM matching process outlined
238 in Section 1: It provides reliable information about the actual presence of an
239 element in reality and thus also on the captured images.

240 Consequently, the proposed method for automated labeling of construc-
241 tion elements uses the data of previously monitored construction sites to-
242 gether with the results from the as-planned vs. as-built comparison to gener-
243 ate valid data sets for the training of neural networks.

244 The proposed workflow is also depicted in Figure 3.

245 As soon as the training is successfully completed, these networks can be
246 used on any construction site for an image based detection of elements.

247 The following subsections describe the process and mathematical back-
248 ground for the projection of construction elements into pictures and the la-
249 beling procedure using these results.

250 *4.1. Camera positions*

251 In the proposed method, the point cloud is produced using photogram-
252 metric methods. In this process, pictures are taken, for example by UAVs
253 (Unmanned aerial vehicles) from different points of view. These pictures can
254 then be used to generate a 3D point cloud if all elements are visible from a

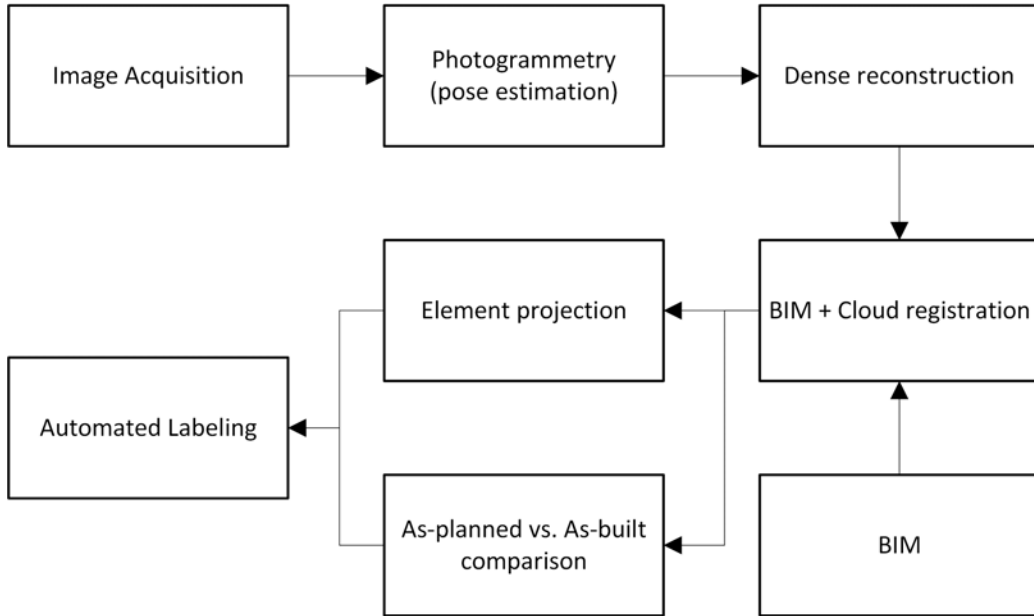


Figure 3: Proposed workflow for the automated labeling toolchain

255 sufficient amount of viewpoints. During the reconstruction process, the camera
 256 positions around the construction site are estimated. This is illustrated
 257 in Fig. 4. This estimation is refined during the dense reconstruction and can
 258 get more accurate by using geodetic reference points on site.

259 4.2. 4D process data and as-planned vs. as-built comparison

260 Building information modeling can be used to combine the geometry of
 261 construction elements with semantic data such as material information but
 262 also process schedules. In the scope of this research, the corresponding process
 263 schedule is connected to all elements, resulting in a fine-grained 4D-BIM
 264 model. This allows identifying all elements that are expected to be built at
 265 each observation time.

266 As depicted in Fig. 5, the software tool used in this research is capable of
 267 integrating the building information model with process data and construction
 268 elements such as scaffolding and formwork.

269 This data is required to define the sets of elements that are used for the
 270 labeling method described in this paper. Since the process schedule may
 271 change during construction, it is crucial to update the schedule permanently
 272 based on the gathered observation data. Since the as-planned vs. as-built

273 comparison has already been conducted for the construction sites in this
 274 research, the results are available for all construction elements. This infor-
 275 mation is crucial since the labeling of elements that were not built yet would
 276 lead to incorrect labels.

277 4.3. Projection

278 Based on the gathered information, it is possible to do a visibility de-
 279 tection by using the camera positions as the point of view, and the process
 280 information to define the set of construction elements, that are meant to
 281 be built. To achieve this, the building model coordinate system needs to
 282 be transformed into the camera coordinate system or vice versa. Several
 283 parameters are needed for this transformation.

284 On the one hand, the intrinsic camera matrix for the distorted images that
 285 projects 3D points in the camera coordinate frame to 2D pixel coordinates
 286 using the focal lengths (F_x, F_y) and the principal point (x_0, y_0) is required.
 287 Additionally, the skew coefficient s_k for the camera is required. This scalar
 288 parameter defines the relation between x and y axis. It is zero if the image
 289 axes are perpendicular. The matrix K can be described as defined in equation
 290 1.

$$K = \begin{bmatrix} F_x & s_k & x_0 \\ 0 & F_y & y_0 \\ 0 & 0 & 1 \end{bmatrix} \quad (1)$$

291 The translation of the camera is defined as:

$$T = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} \quad (2)$$

292 Additionally, the rotation matrix for each image as defined in equation 3
 293 is needed.

$$R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \quad (3)$$

294 Both, translation and rotation can be described in one 3 x 4 matrix:

$$RT = \begin{bmatrix} r_{11} & r_{12} & r_{13} & T_1 \\ r_{21} & r_{22} & r_{23} & T_2 \\ r_{31} & r_{32} & r_{33} & T_3 \end{bmatrix} \quad (4)$$

295 Using the model coordinates of all triangulated construction elements,
 296 it is possible to calculate the projection of each element into the camera
 297 coordinate system and therefore overlay the model projection and the corre-
 298 sponding picture taken from the point of observation with equation 5.

$$t = K * RT * p; \quad (5)$$

The resulting 2D coordinates that are rendered into the picture are calcu-
 lated by using the vector t and getting the x and y coordinates by calculating

$$x = t[0]/t[2] \quad (6)$$

and

$$y = t[1]/t[2] \quad (7)$$

299 This is done for each point belonging to the triangulated geometry rep-
 300 resentation of all construction elements.

301 As visible in Fig. 6 for an analytical column, the projection works as
 302 expected and helps to identify the respective construction element in the
 303 recorded picture. The mentioned calculations need to include an optional
 304 transformation and rotation if the model is geo-referenced and thus the two
 305 coordinate systems differ broadly.

306 4.4. Render model based on camera position

307 The algorithm introduced in section 4.3 enables the element-wise render-
 308 ing of all construction elements in the respective coordinate system. To get
 309 a rendered image of all visible construction elements, the following steps are
 310 carried out:

311 While all geometric information is available, three problems need to be
 312 solved for an accurate rendering of all construction elements:

- 313 1. For triangulated elements, only the boundaries are known. However,
 314 the whole surface needs to be rendered correctly.
- 315 2. The rendered surface needs to be connected to the corresponding ele-
 316 ment since this information is crucial for a proper visibility analysis

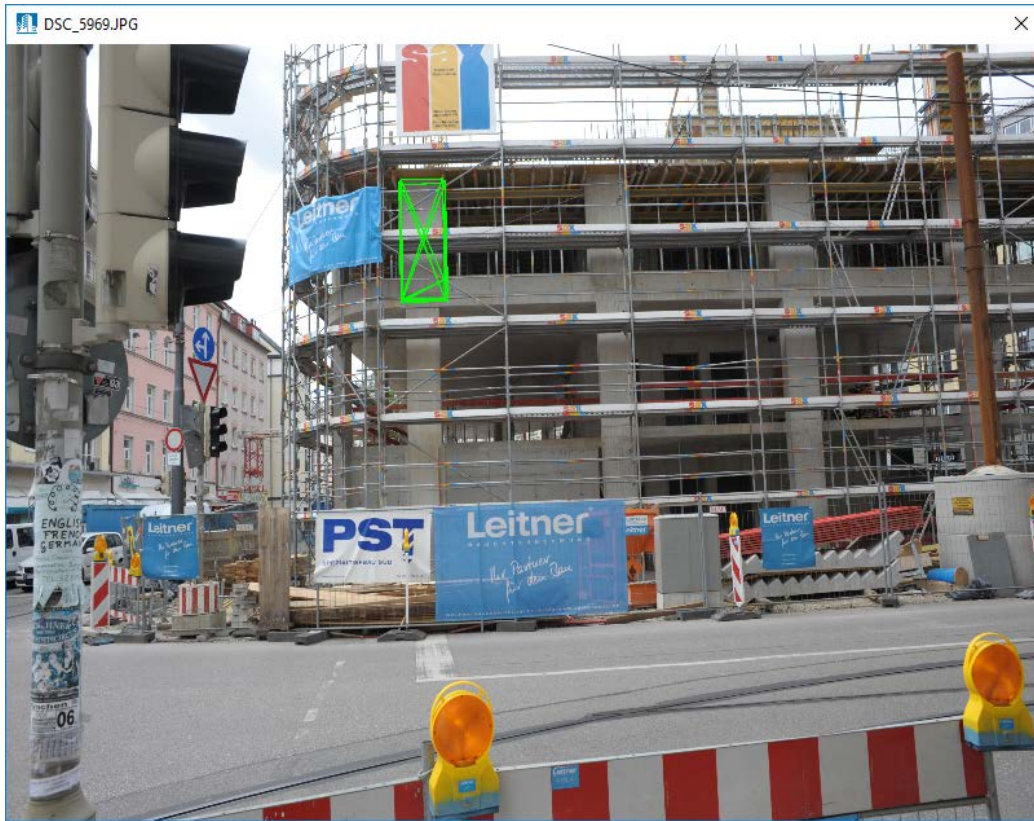


Figure 6: Sample of projected, triangulated column geometry into a corresponding picture

Algorithm 1 Pseudo code for rendering an image of all visible elements

```
1: procedure RENDERVISIBLEELEMENTS
2:    $O \leftarrow$  set of all observations of the construction site
3:    $I \leftarrow$  set of all images of current observation
4:    $E \leftarrow$  set of all construction elements
5:    $C \leftarrow$  set of all coordinates of the triangulated surfaces
6:    $d \leftarrow$  distance of element to corresponding camera position
7:   for all  $O$  do
8:     for all  $I$  do
9:       for all  $E$  do
10:        for all  $C$  do
11:          if isvisible(c) then  $P(x,y,d,color) = \text{projection}(c)$ ;
12:        for all Pixels do
            $p_{min} = \min(P(d))$ ;
            $p(x,y) = p(color)$ ;
```

317 3. Elements may blend over from the viewpoint in some circumstances.
318 This needs to be addressed to get a correct rendering.

319 The first issue is solved by applying necessary inside/outside tests for
320 points inside a bounding box around each triangle. This is combined with
321 min/max tests to verify that all points are inside the given coordinate sys-
322 tem of the current picture. The second issue is addressed by assigning an
323 individual color in the RGB color range to every construction element. This
324 allows identifying each element after the rendering is finished.

325 The third issue is solved by applying the Painter's algorithm [31] to each
326 pixel in the given picture. In the given challenge, the distance to the point of
327 view is stored for the current construction element and the color information
328 is replaced in case an element has a smaller distance to the point of view and
329 is also visible in the same pixel of the picture.

330 The applied algorithms result in a rendering as seen in Fig. 7.

331 After applying this technique to all observations and all camera positions,
332 a distinct list of all visible construction elements can be generated by iterating
333 over all pixels of each rendered image. The color of each pixel is assigned to
334 a construction element, and since the painters' algorithm is applied, only the
335 element is visible, that has the lowest distance to the point of observation.
336 Therefore, all visible, non-occluded elements can be determined with this

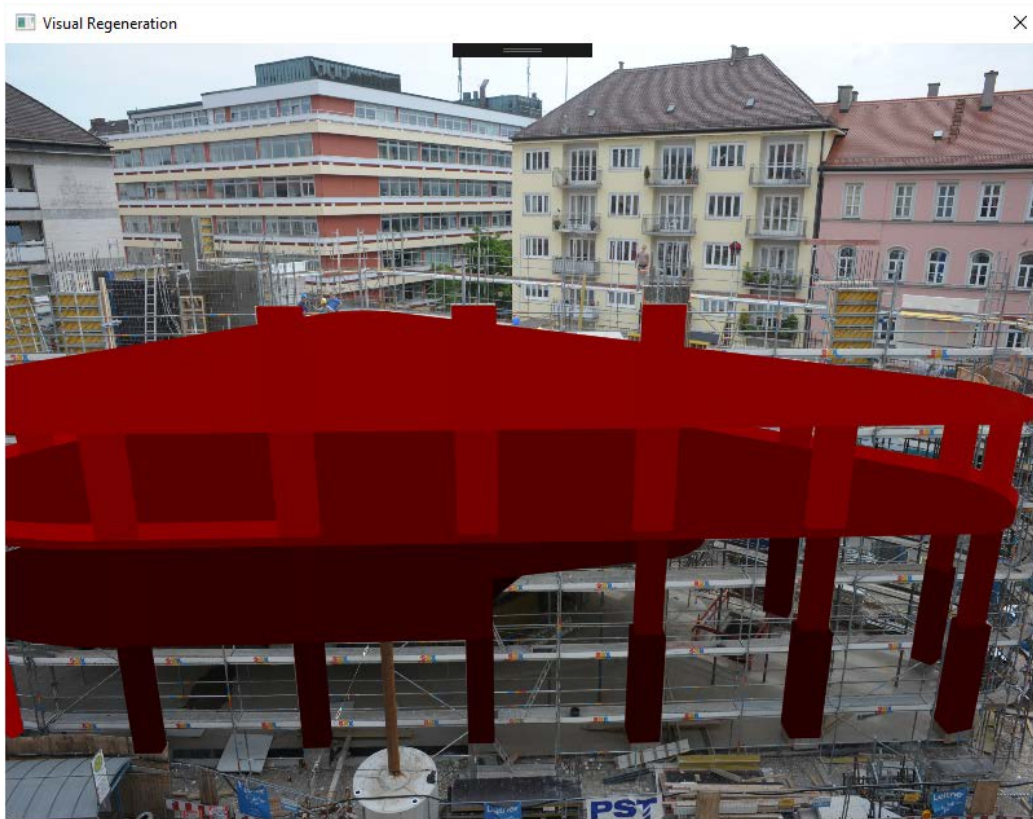


Figure 7: Using projection methodology for model rendering based on the Painter's Algorithm and 4D semantic information

337 method.

338 *4.5. Generating Labels for Machine Learning*

339 Since machine learning tasks require large training sets for the learning
340 procedure, the labeling and pre-processing of suitable data play a crucial
341 role.

342 Labeling for ML depends on the desired output of the ML system. A basic
343 ML system for classification is only capable of making general statements on
344 the content of an image and hence only requires a set of images containing
345 the classification category as training input. On the other hand, a system for
346 semantic segmentation can predict the exact location and also the amount of
347 (multiple) elements in one image. Labeling for this class of systems requires
348 detailed convex hull polygons around all instances of elements. Additionally,
349 the category for each label needs to be defined.

350 The automated labeling process presented here builds on the previously
351 presented projection algorithm and is capable of generating labels for all
352 sorts of ML systems starting from necessary bounding boxes up to detailed
353 convex hulls around individual element instances. Image-based labeling is
354 realized by defining a polygon line around each object and associating a
355 corresponding category with this label. The polygon label can be generated
356 by the above-mentioned projection and fits precisely around the shape of
357 each construction element. The defining element category can be extracted
358 from the semantic information provided by the building information model.
359 Since geometry and semantic data are connected, any additional information
360 can be added to the generated labels.

361 Besides using the mathematical algorithms for projection, also the results
362 of the visibility analysis are essential. As discussed before and depicted in
363 Figure 8, labeling cannot only rely on all available elements. A prominent but
364 noteworthy factor is the actual presence of the labeled element. The element
365 must have been built to generate an image valid for training or testing. By
366 extracting this information from the as-planned vs. as-built comparison, the
367 set of available elements is reduced to the set of detected elements. In the
368 next step, the set needs to be further reduced to the set of visible elements
369 for each picture.

370 To sum up the labeling process, the following method is proposed:

371 The proposed method works for all kinds of label requirements. To illus-
372 trate this, Table 1 shows sample labels for a Classification Network (which

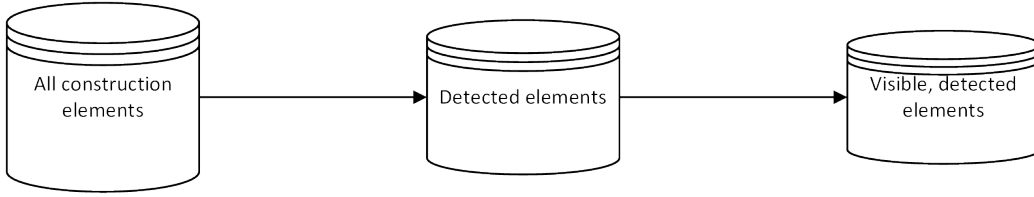


Figure 8: Considered set of elements based on previous results from as-planned vs. as-built comparison. For CV-based methods, visibility plays a crucial role, leading to reduced data sets from construction monitoring.

Algorithm 2 Pseudo code for labeling all visible elements in an image

```

1: procedure LABELVISIBLEELEMENTS
2:    $I \leftarrow$  set of all images
3:    $List \langle element, List \langle P \rangle \rangle LabelList \leftarrow$  set of all labels
4:   for all I do
5:      $E \leftarrow$  set of all construction elements, visible in current picture
6:     for all E do
7:        $List \langle P \rangle ConvexHull = GetConvexHull();$ 
        $LabelList.Add(E.elementtype, ConvexHull);$ 
  
```

373 usually requires image snippets with bounding boxes) and semantic segmen-
 374 tation (which usually requires polygon lines and the corresponding images).
 375 The sample shown here uses the well known COCO format. For better un-
 376 derstanding, a graphic labeled image for semantic segmentation is added,
 377 too.

378 Current best practice in machine learning proposes to split the labeled
 379 data-set into a set of training data for the actual training process, a set of
 380 validation that does not contain any data from the training set to validate the
 381 current training rates. Finally, a set for testing that is not used for training
 382 or validation at all is used for checking the overall performance of the neural
 383 network without further changing the learning parameters.

384 Hence, the labeled images are split randomly into the mentioned cate-
 385 gories to fulfill this requirement.

386 5. Case study

387 To prove the introduced methods, the following case studies were con-
 388 ducted:





Category	Classification	semantic seg. [JSON coords]	seg. image
column		[[3778, 1230, 3810, 1230, 3834, 1231, 3837, 1230, 3840, 1230, 3854, 983, 3852, 984, 3848, 984, 3791, 985]]	
formwork		[[2662, 1662, 2666, 1682, 2703, 1682, 2704, 1323, 2702, 1319, 2699, 1314, 2663, 1314]]	

Table 1: Sample labels of two categories for different ML use cases

389 5.1. BIM element projection onto images

390 The developed methodology has been applied to several construction
 391 sites.

392 As depicted in Fig. 2, most observations lack details at some point and
 393 have mostly occluded areas due to the observation methods. In very dis-
 394 advantageous observations, the detection rate can drop down to 50% of the
 395 overall built construction elements. In this case, the detection rate d describes
 396 the percentage of elements that the proposed method marked as detected over
 397 the ground truth of all elements that were built. The latter set of elements
 398 has been acquired manually in order to verify used algorithms. With the
 399 help of the presented methods, these rates can be explained since most of
 400 the undetected elements were not visible from the observation points. To
 401 quantify the efficiency of an algorithm for as-planned vs. as-built detection,
 402 it is essential to have a valid ground truth to allow an unbiased evaluation
 403 of the used methods. This approach helps to quantify the used methods
 404 correctly.

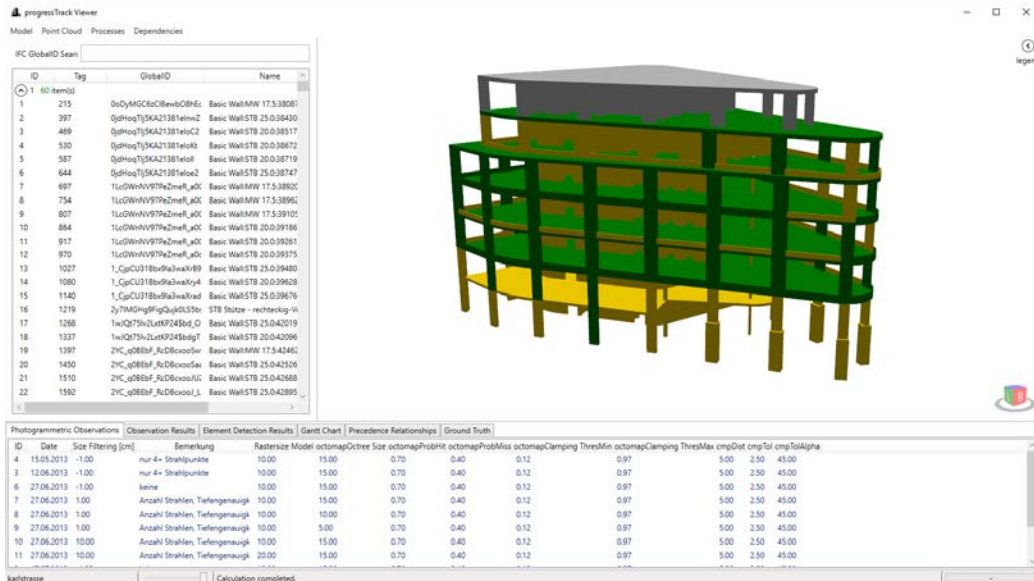


Figure 9: Detected construction elements from one observation. Green elements were successfully detected, yellow elements were not detected but are built.

405 This concept is illustrated in Fig. 9. The green elements were detected
 406 correctly. The yellow elements, however, are built but were not detected.

407 This is because the inner walls were not visible from a sufficient number of
408 viewpoints. Thus there were not enough points in the corresponding point
409 cloud that allowed to validate the existence of the elements. However, these
410 elements were identified as not visible using the method introduced earlier
411 in this paper.

412 *5.2. Automated labeling and validation*

413 After successfully testing the projection, the actual labeling is performed
414 being the key contribution of this paper.

415 Many currently used CNNs rely on the COCO Data-set [32]. Facebook’s
416 Mask R-CNN [33] has provided promising results for machine learning in pre-
417 vious applications. The network itself also relies on the COCO data format
418 as a basis. Thus, the authors chose this schema as a basis for the generation
419 of the labels. This schema requires a defined structure for all labels, including
420 information about each image (id, width, height, license, date captured), all
421 annotations (id, corresponding image, label category, label polygon, bound-
422 ing box, ...) as well as the defined categories (in this case for example walls
423 or columns, all represented by individual IDs).

424 The construction projects on which the developed methods have been
425 applied involve mainly the production of concrete elements. The following
426 construction elements and temporary elements were modeled in the corre-
427 sponding BIM:

- 428 • columns
- 429 • walls
- 430 • formworks
- 431 • slabs
- 432 • roofs
- 433 • stairs

434 The proposed methods were tested on observation data from multiple
435 construction sites, resulting in 32,787 labeled construction elements on 1,300
436 images. The machine used for this test is a Windows 10 system equipped
437 with an Intel Xeon E5-1630 CPU @ 3.70GHz, 16 GB DDR4-RAM, AMD
438 Fire Pro W4100 2GB RAM, and a 10Gbit network connection. The entire

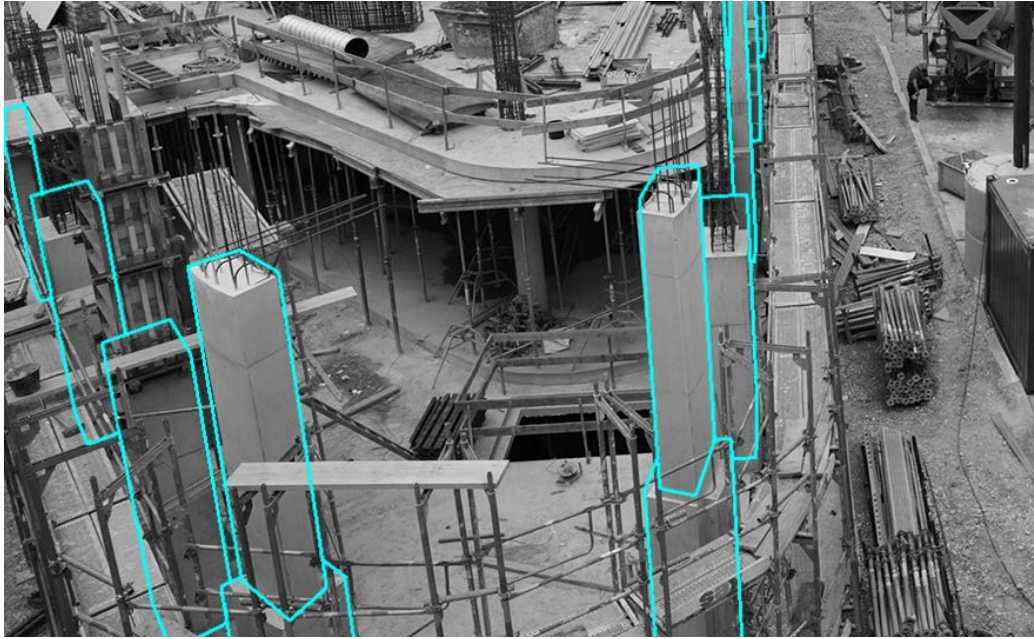


Figure 10: Sample sub-set of auto-labeled columns in one picture from a construction site.



Figure 11: Sample sub-set of auto-labeled columns and walls in another picture from a construction site.

439 automated labeling process took around 20 minutes, outperforming manual
440 labeling significantly. During this time, all images (close to 9 GB) were
441 downloaded from a NAS (Network attached storage) and randomly added to
442 the training, validation, and test data sets. Additionally, the corresponding
443 label files in JSON format were generated. A sample visualization of the
444 generated labels for one picture is depicted in Figure 10 and 11, showing
445 exported columns with their respective convex hull label around them.

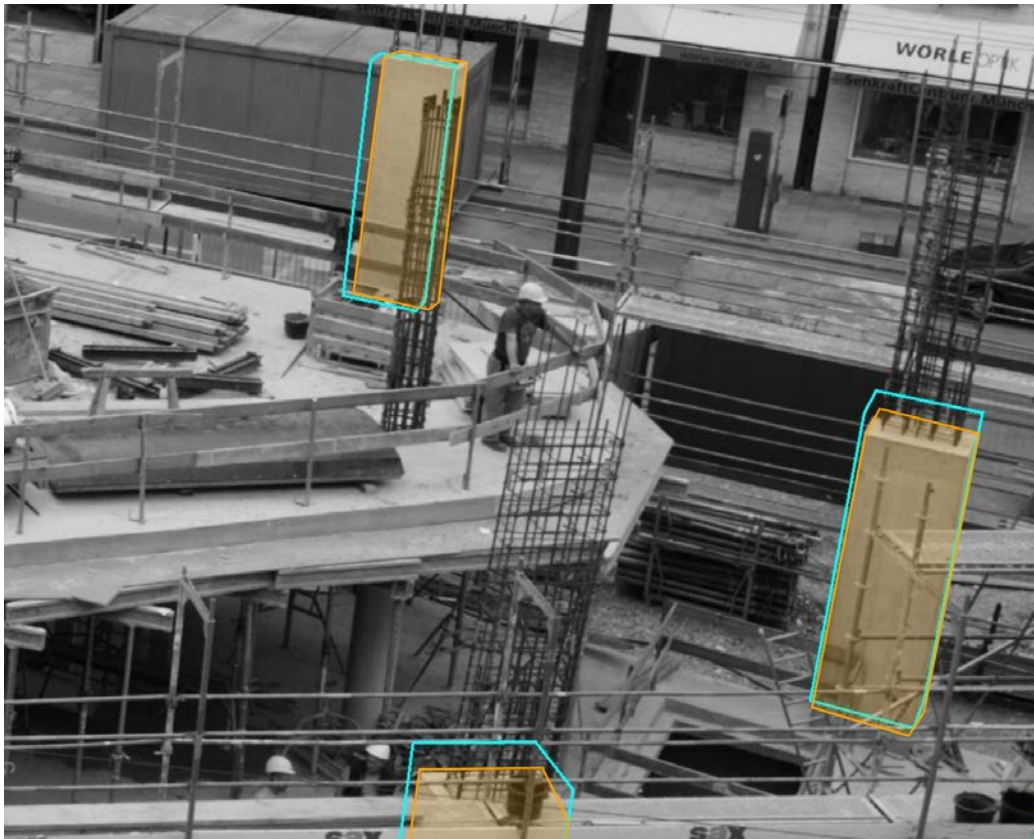


Figure 12: Validation of label correctness with cyan poly-lines representing the automatically generated labels and orange poly-lines representing the manually generated validation set.

446 The method was validated through human evaluation on all labels for the
447 tested construction sites. The label projection worked without failure for all
448 built construction elements in terms of generating a valid convex hull as the
449 existing elements have been verified against a manually created ground truth.

450 Since no issues were found in a set of over 32,000 snippets, the projection
451 can be regarded as working correctly. However, as depicted in Fig. 12, the
452 automatically generated labels (cyan poly-lines) have a slight deviation from
453 the actual construction elements.

454 This deviation can have multiple reasons:

- 455 • errors in pose estimation during Structure-from-Motion
- 456 • large scale deviations when using real-world coordinates
- 457 • construction inaccuracies
- 458 • modeling inaccuracies

459 Since all elements were validated in the as-planned vs. as-built compari-
460 son, allowing for only very minor construction inaccuracies, construction in-
461 accuracies can be disregarded in this research. Otherwise, the element would
462 not have been classified as "built" and would not have been labeled at all.
463 Thus, the deviations, in this case, are minor and arise from an aggregation
464 of the mentioned reasons. To quantify the introduced error, a set of 1.000
465 elements have been labeled manually and tested against the automatically
466 created labels. The labels were then compared pixel-wise. The overall accu-
467 racy of the automated system was measured by calculating the overlapping
468 area I of the resulting labels of both labeling methods over the manually
469 labeled area:

$$p_o = I/A_{manuallabel} \quad (8)$$

470 with

$$I = A_{autolabel} \cap A_{manuallabel} \quad (9)$$

471 The resulting accuracy p_o had an average of 91.7% overlap over all checked
472 labels, constantly lying within the bounds of 85% and 97%. The overlap rates
473 give promising results and make the labels usable for machine learning tasks.
474 Rates could be further improved by taking more pictures for the Structure-
475 from-Motion process and enhancing the resulting camera pose estimation.

476 6. Discussion

477 For improving the reliability of construction progress monitoring, this
478 paper introduces a novel concept for automating the labeling process of con-
479 struction site images. It is based on fusing information available from the
480 photogrammetric process (images and relative position of the camera) and
481 the information available from the 4D BIM (object type, object position).
482 Since the BIM and the resulting point cloud are aligned, a digital element
483 can be projected onto the image, initially taken for the photogrammetric
484 process. Also, matching the point cloud and the BIM allows to make sure
485 that only images are considered where the elements under consideration exist
486 in reality.

487 From the projected BIM elements, it is possible to automatically con-
488 nect the covered image segments with the semantic information provided
489 by the Building Information Model. Since the introduced as-planned vs.
490 as-built comparison also offers valuable information on the presence of all
491 elements, the labels can be further refined regarding possible occlusions. As
492 a valid label should only be applied to an at least partially visible element,
493 the gathered knowledge from the previously applied as-planned vs. as-built
494 comparison makes this automated approach even more accurate. Since the
495 comparisons' resulting elements are built at the correct positions, the labels
496 are also correct. On the downside, only elements that were built as-planned
497 can be labeled.

498 The sample-based validation showed over 91% pixel-wise accuracy of the
499 automated procedure when tested against manual labeling procedures. A
500 previously tested, manual labeling approach took over 100 working hours to
501 accurately label only one category of elements. Labeling and generating the
502 corresponding images folders for this case study took around 20 minutes,
503 including downloading of 9 GB of pictures from a remote NAS folder which
504 takes over 90% of the time used. Additionally, several studies show that
505 manual labeling is also introducing a range of errors due to missed elements
506 or inaccurately labeled elements. As the correct identification of construction
507 elements also requires technical personnel [34], labeling is hugely cost inten-
508 sive and danger of bore-out to this group of workers due to the repetitive
509 work.

510 The construction sites used for this process are located in Germany and
511 apply in-situ concrete pouring as the primary construction methodology.
512 Consequently, the resulting labels and especially the trained network, will

513 only be able to detect construction elements from this domain of manufac-
514 turing. However, the presented approach can be easily extended by also
515 including construction sites from other countries or other construction tech-
516 niques.

517 Future steps of this research will focus on creating a CNN for detecting
518 the most important construction elements on construction sites. The final
519 objective is to enable a utterly image-based construction monitoring process
520 in the future.

521 **Acknowledgments**

522 We thank the Leibniz Supercomputing Centre (LRZ) of the Bavarian
523 Academy of Sciences and Humanities (BAdW) for the support and provi-
524 sion of computing infrastructure essential to this publication. We thank the
525 German Research Foundation (DFG) for funding the initial phases of the
526 project.

527 Additionally, we would like to thank all construction companies, pro-
528 viding data and/or support to the research conducted in this field, includ-
529 ing: Leitner GmbH & Co Bauunternehmung KG, Kuehn Malvezzi Archi-
530 tects, Staatliches Bauamt Muenchen, Baugesellschaft Brunner+Co, BKL
531 (Baukranlogistik GmbH), Geiger Gruppe, Baureferat H5, Landeshauptstadt
532 Muenchen, Baugesellschaft Mickan mbH & Co KG, h4a Architekten, Wenzel
533 + Wenzel and Stadtvermessungsamt Muenchen.

534 **References**

- 535 [1] A. Braun, S. Tuttas, A. Borrmann, U. Stilla, Automated Progress Mon-
536 itoring Based on Photogrammetric Point Clouds and Precedence Rela-
537 tionship Graphs, Proceedings of the 32nd International Symposium on
538 Automation and Robotics in Construction and Mining (ISARC 2015)
539 (2017) 274–280. doi:10.22260/isarc2015/0034.
- 540 [2] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich Feature Hierarchies
541 for Accurate Object Detection and Semantic Segmentation, in: 2014
542 IEEE Conference on Computer Vision and Pattern Recognition, IEEE,
543 2014, pp. 580–587. URL: [http://ieeexplore.ieee.org/document/
544 6909475/](http://ieeexplore.ieee.org/document/6909475/). doi:10.1109/CVPR.2014.81.

- 545 [3] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Im-
546 age Recognition, in: 2016 IEEE Conference on Computer Vision and
547 Pattern Recognition (CVPR), IEEE, 2016, pp. 770–778. URL: [http://](http://ieeexplore.ieee.org/document/7780459/)
548 ieeexplore.ieee.org/document/7780459/. doi:10.1109/CVPR.2016.
549 90.
- 550 [4] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You Only Look
551 Once: Unified, Real-Time Object Detection, in: 2016 IEEE Confer-
552 ence on Computer Vision and Pattern Recognition (CVPR), IEEE,
553 2016, pp. 779–788. URL: [http://ieeexplore.ieee.org/document/](http://ieeexplore.ieee.org/document/7780460/)
554 [7780460/](http://ieeexplore.ieee.org/document/7780460/). doi:10.1109/CVPR.2016.91.
- 555 [5] A. Dimitrov, M. Golparvar-Fard, Vision-based material recognition for
556 automated monitoring of construction progress and generating building
557 information modeling from unordered site image collections, *Advanced*
558 *Engineering Informatics* 28 (2014) 37–49. doi:10.1016/j.aei.2013.11.
559 002.
- 560 [6] I. Brilakis, L. Soibelman, Y. Shinagawa, Material-Based Construction
561 Site Image Retrieval, *Journal of Computing in Civil Engineering* 19
562 (2005) 341–355. doi:10.1061/(asce)0887-3801(2005)19:4(341).
- 563 [7] J. J. Lin, K. K. Han, M. Golparvar-Fard, A Framework for Model-Driven
564 Acquisition and Analytics of Visual Data Using UAVs for Automated
565 Construction Progress Monitoring, in: *Computing in Civil Engineering*
566 2015, American Society of Civil Engineers, Reston, VA, 2015, pp. 156–
567 164. URL: [http://ascelibrary.org/doi/10.1061/9780784479247.](http://ascelibrary.org/doi/10.1061/9780784479247.020)
568 020. doi:10.1061/9780784479247.020.
- 569 [8] T. Omar, M. L. Nehdi, Data acquisition technologies for construc-
570 tion progress tracking, *Automation in Construction* 70 (2016) 143–155.
571 doi:10.1016/j.autcon.2016.06.016.
- 572 [9] F. Bosche, C. T. Haas, Automated retrieval of 3D CAD model objects
573 in construction range images, *Automation in Construction* 17 (2008)
574 499–512. doi:10.1016/j.autcon.2007.09.001.
- 575 [10] F. Bosché, Plane-based registration of construction laser scans with
576 3D/4D building models, *Advanced Engineering Informatics* 26 (2012)
577 90–102. doi:10.1016/j.aei.2011.08.009.

- 578 [11] Y. Turkan, F. Bosché, C. T. Haas, R. Haas, Automated progress track-
579 ing using 4D schedule and 3D sensing technologies, *Automation in Con-*
580 *struction* 22 (2012) 414–421. doi:10.1016/j.autcon.2011.10.003.
- 581 [12] C. C. Kim, H. Son, C. C. Kim, Fully automated registration of 3D data
582 to a 3D CAD model for project progress monitoring, *Automation in*
583 *Construction* 35 (2013) 587–594. doi:10.1016/j.autcon.2013.01.005.
- 584 [13] C. Zhang, D. Arditi, Automated progress control using laser scanning
585 technology, *Automation in Construction* 36 (2013) 108–116. doi:10.
586 1016/j.autcon.2013.08.012.
- 587 [14] Y. Ibrahim, T. Lukins, X. Zhang, E. Trucco, a. Kaka, Towards auto-
588 mated progress assessment of workpackage components in construction
589 projects using computer vision, *Advanced Engineering Informatics* 23
590 (2009) 93–103. doi:10.1016/j.aei.2008.07.002.
- 591 [15] C. C. Kim, H. Son, C. C. Kim, Automated construction progress mea-
592 surement using a 4D building information model and 3D data, *Automa-*
593 *tion in Construction* 31 (2013) 75–82. doi:10.1016/j.autcon.2012.11.
594 041.
- 595 [16] H. Son, C. Kim, 3D structural component recognition and modeling
596 method using color and 3D data for construction progress monitor-
597 ing, *Automation in Construction* 19 (2010) 844–854. doi:10.1016/j.
598 autcon.2010.03.003.
- 599 [17] I. Brilakis, H. Fathi, A. Rashidi, Progressive 3D reconstruction of infras-
600 tructure with videogrammetry, *Automation in Construction* 20 (2011)
601 884–895. doi:10.1016/j.autcon.2011.03.005.
- 602 [18] I. Brilakis, S. German, Z. Zhu, Visual Pattern Recognition Models
603 for Remote Sensing of Civil Infrastructure, *Journal of Computing*
604 *in Civil Engineering* 25 (2011) 388–393. doi:10.1061/(ASCE)CP.1943-
605 5487.0000104.
- 606 [19] A. Rashidi, I. Brilakis, P. Vela, Generating Absolute-Scale Point Cloud
607 Data of Built Infrastructure Scenes Using a Monocular Camera Setting,
608 *Journal of Computing in Civil Engineering* 29 (2015) 04014089. doi:10.
609 1061/(ASCE)CP.1943-5487.0000414.

- 610 [20] M. Golparvar-Fard, M. Asce, F. Peña-Mora, S. Savarese, M. Asce,
611 F. Peña-Mora, S. Savarese, Integrated Sequential As-Built and As-
612 Planned Representation with Tools in Support of Decision-Making
613 Tasks in the AEC/FM Industry, *Journal of Construction Engineering*
614 *and Management* 137 (2011) 1099–1116. doi:10.1061/(ASCE)CO.1943-
615 7862.0000371.
- 616 [21] M. Golparvar-Fard, F. Pena-Mora, S. Savarese, Monitoring changes
617 of 3D building elements from unordered photo collections, *Computer*
618 *Vision Workshops (ICCV Workshops)*, 2011 IEEE International Con-
619 ference on (2011) 249–256. doi:10.1109/ICCVW.2011.6130250.
- 620 [22] M. Golparvar-Fard, F. Peña-Mora, S. Savarese, Automated Progress
621 Monitoring Using Unordered Daily Construction Photographs and IFC-
622 Based Building Information Models, *Journal of Computing in Civil*
623 *Engineering* 29 (2015) 04014025. doi:10.1061/(ASCE)CP.1943-5487.
624 0000205.
- 625 [23] K. Han, M. Golparvar-Fard, Crowdsourcing BIM-guided collection of
626 construction material library from site photologs, *Visualization in En-*
627 *gineering* 5 (2017) 14. doi:10.1186/s40327-017-0052-3.
- 628 [24] S. Chi, C. H. Caldas, Automated Object Identification Using Opti-
629 cal Video Cameras on Construction Sites, *Computer-Aided Civil and*
630 *Infrastructure Engineering* 26 (2011) 368–380. doi:10.1111/j.1467-
631 8667.2010.00690.x.
- 632 [25] C. Kropp, C. Koch, M. König, Interior construction state recognition
633 with 4D BIM registered image sequences, *Automation in Construction*
634 86 (2018) 11–32. doi:10.1016/j.autcon.2017.10.027.
- 635 [26] C. Kim, B. Kim, H. Kim, 4D CAD model updating using image
636 processing-based construction progress monitoring, *Automation in Con-*
637 *struction* 35 (2013) 44–52. doi:10.1016/j.autcon.2013.03.005.
- 638 [27] B. Akinci, F. Boukamp, C. Gordon, D. Huber, C. Lyons, K. Park, A
639 formalism for utilization of sensor systems and integrated project models
640 for active construction quality control, *Automation in Construction* 15
641 (2006) 124–138. doi:10.1016/j.autcon.2005.01.008.

- 642 [28] C. Kropp, C. Koch, M. König, Drywall State Detection in Image Data
643 for Automatic Indoor Progress Monitoring, in: Computing in Civil and
644 Building Engineering (2014), November 2015, American Society of Civil
645 Engineers, Reston, VA, 2014, pp. 347–354. URL: <http://ascelibrary.org/doi/10.1061/9780784413616.044>. doi:10.1061/9780784413616.044.
647
- 648 [29] A. Braun, S. Tuttas, U. Stilla, A. Borrmann, Process- and Computer
649 Vision-based Detection of As-Built Components on Construction Sites,
650 in: 2018 Proceedings of the 35th ISARC, Berlin, Germany, 2018, p. 7.
651 doi:10.22260/ISARC2018/0091.
- 652 [30] C. Wu, Towards linear-time incremental structure from motion,
653 in: Proceedings - 2013 International Conference on 3D Vision, 3DV
654 2013, IEEE, 2013, pp. 127–134. URL: <http://ieeexplore.ieee.org/articleDetails.jsp?arnumber=6599068>. doi:10.1109/3DV.2013.25.
655
- 656 [31] T. T. Elvins, A survey of algorithms for volume visualization, ACM SIG-
657 GRAPH Computer Graphics 26 (2005) 194–201. doi:10.1145/142413.142427.
658
- 659 [32] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays,
660 P. Perona, D. Ramanan, C. L. Zitnick, P. Dollár, Microsoft COCO:
661 Common Objects in Context, Proceedings of the IEEE Computer So-
662 ciety Conference on Computer Vision and Pattern Recognition (2014)
663 3686–3693. doi:10.1109/CVPR.2014.471.
- 664 [33] K. He, G. Gkioxari, P. Dollar, R. Girshick, Mask R-CNN, Proceedings
665 of the IEEE International Conference on Computer Vision 2017-Octob-
666 (2017) 2980–2988. doi:10.1109/ICCV.2017.322.
- 667 [34] K. K. Han, M. Golparvar-Fard, Potential of big visual data and build-
668 ing information modeling for construction performance analytics: An
669 exploratory study, Automation in Construction 73 (2017) 184–198.
670 doi:10.1016/j.autcon.2016.11.004.