

Estimating Energy Consumption of Residential Buildings at Scale with Drive-by Image Capture

WOC Ward X Li Y Sun M Dai H Arbabi D Densley Tingley
M Mayfield
Department of Civil and Structural Engineering, The University of Sheffield, UK

Abstract

Data driven approaches to addressing climate change are increasingly becoming a necessary solution to deal with the scope and scale of interventions required to reach net zero. In the UK, housing contributes to over 30% of the national energy consumption, and a massive rollout of retrofit is needed to meet government targets for net zero by 2050. This paper introduces a modular framework for quantifying building features using drive-by image capture and utilising them to estimate energy consumption. The framework is demonstrated on a case study of houses in a UK neighbourhood, showing that it can perform comparatively with gold standard datasets. The paper reflects on the modularity of the proposed framework, potential extensions and applications, and highlights the need for robust data collection in the pursuit of efficient, large-scale interventions.

1 Introduction

Buildings and their operation contribute to nearly 17% of global carbon emissions (UNEP, 2020). Of these emissions, 61% can be mitigated, according to the IPCC (2022), with the largest share of mitigation potential coming from the retrofit of existing buildings in developed countries such as the UK. The report also highlights that the next decade is critical for building technical capacity to ensure this potential is realised (IPPC, 2022). Technical solutions to the mass implementation of retrofit require robust, large scale data and modelling.

Large scale modelling of residential buildings with a resolution of information at the individual building level requires high volumes of data. Capturing and processing high quality data that can be used in decision making both reliably and efficiently, in high volumes, will require a substantial degree of automation. However, access to reliable sources of built environment data can be a challenge. Building stock models, for example, have been developed for use in modelling energy usage and occupant behaviour at an individual building level, however such methods have relied on a set of predefined archetypes (Shorrocks et al., 2005). Such archetypes, that describe e.g. age cohorts, can miss particular nuances in different construction types, or building performance. In Great Britain (GB), datasets such as those provided by Ordnance Survey (Ordnance Survey, 2022b) and Verisk (Verisk, 2022) provide attributes for individual properties on a national scale, including building footprints, building heights and usage. Previous work has looked at this data, along with aerial point cloud data, to produce city-level stock models (Steadman et al., 2020),

31 however while the aerial data can provide large scale topographic information, there are limited
32 resources to provide information on facades and other street-level features that are essential in
33 aiding understanding of the urban environment, something that is essential for enacting retrofit
34 measures.

35 Energy consumption per unit floor area is a metric used to assess the energy efficiency of
36 a building. In the UK, this is most commonly reported in the form of energy performance
37 certificates (EPC). The generation of EPCs requires a manual survey of the property, which
38 is then used to input information into the so-called Standard Assessment Procedure (SAP) to
39 estimate energy consumption (BRE, 2012). These assessments take approximately 45 minutes
40 per building to complete and are conducted as required by law: most commonly when a property
41 is sold, or every ten years in the case of rental properties (ReallyMoving.com, 2022). However,
42 reports of issues in EPC reporting are widespread (Hardy and Glew, 2019). For example, the
43 Retrofit Playbook, a guide to retrofit for policy makers in the UK published by the UK Green
44 Buildings Council, describes EPCs as “not fit for purpose”, and highlight this as a barrier to
45 enacting home retrofit in the UK (UKGBC, 2021). To help overcome these barriers, a framework
46 is presented in this paper that uses drive-by captured image data to generate energy models at
47 a high scale.

48 The framework is designed as a set of modular components, defined in terms of their input and
49 output, with the aim to simulate energy consumption for individual properties at scale. The
50 framework is a data-driven approach that utilises computer vision techniques, including machine
51 learning and 3-D reconstruction, to measure and assess buildings using street-level images and
52 use this to build models to simulate energy consumption. Each component of the framework is
53 discussed, and the approach applied to a case study of residential buildings in a neighbourhood
54 in the UK. The capabilities of the proposed framework are compared against available data,
55 and the results critically appraised in terms of performance against this existing data. The
56 framework also utilises a modular approach that could be used in conjunction with other methods.
57 Furthermore, the paper addresses the potential scope for extension and generalisation of the
58 framework with additional modalities of data, such as thermal imaging, and the application to
59 other quantification problems in built environment research.

60 **2 Framework Design**

61 This section presents a detailed overview of the proposed framework for capturing and localising
62 street-level images and using them to estimate the physical properties of houses for simulation of
63 energy consumption. Figure 1 provides an illustration of the framework as a pipeline of modular
64 components mapping from data collection and localisation to energy modelling. Each framework
65 component highlighted in Figure 1 is described in the following sections.

66 **2.1 Data Collection and Localisation**

67 A scalable platform for estimating building properties and energy consumption requires a scalable
68 solution to data generation and processing. Data analysis and decision making that can be

69 performed at neighbourhood- or city-scale requires large quantities of high quality, localised
 70 data.

71 2.1.1 Capture

72 Drive-by data capture is not uncommon, and has been used for applications such as mapping
 73 (Anguelov et al., 2010) or in the development of self-driving vehicle technology (Gwak et al.,
 74 2019). To effectively perform feature and geometry extraction, image data needs to be available
 75 in high temporal and spatial resolution. In this research, a bespoke mobile sensing vehicle is
 76 used to prototype the proposed framework (Arbabi et al., 2021).

77 Image data is captured by driving the sensing vehicle along residential streets. The mobile
 78 sensing vehicle uses a multi-camera rig¹ to capture spherical image data using five radial cameras

¹Teledyne FLIR Ladybug5+

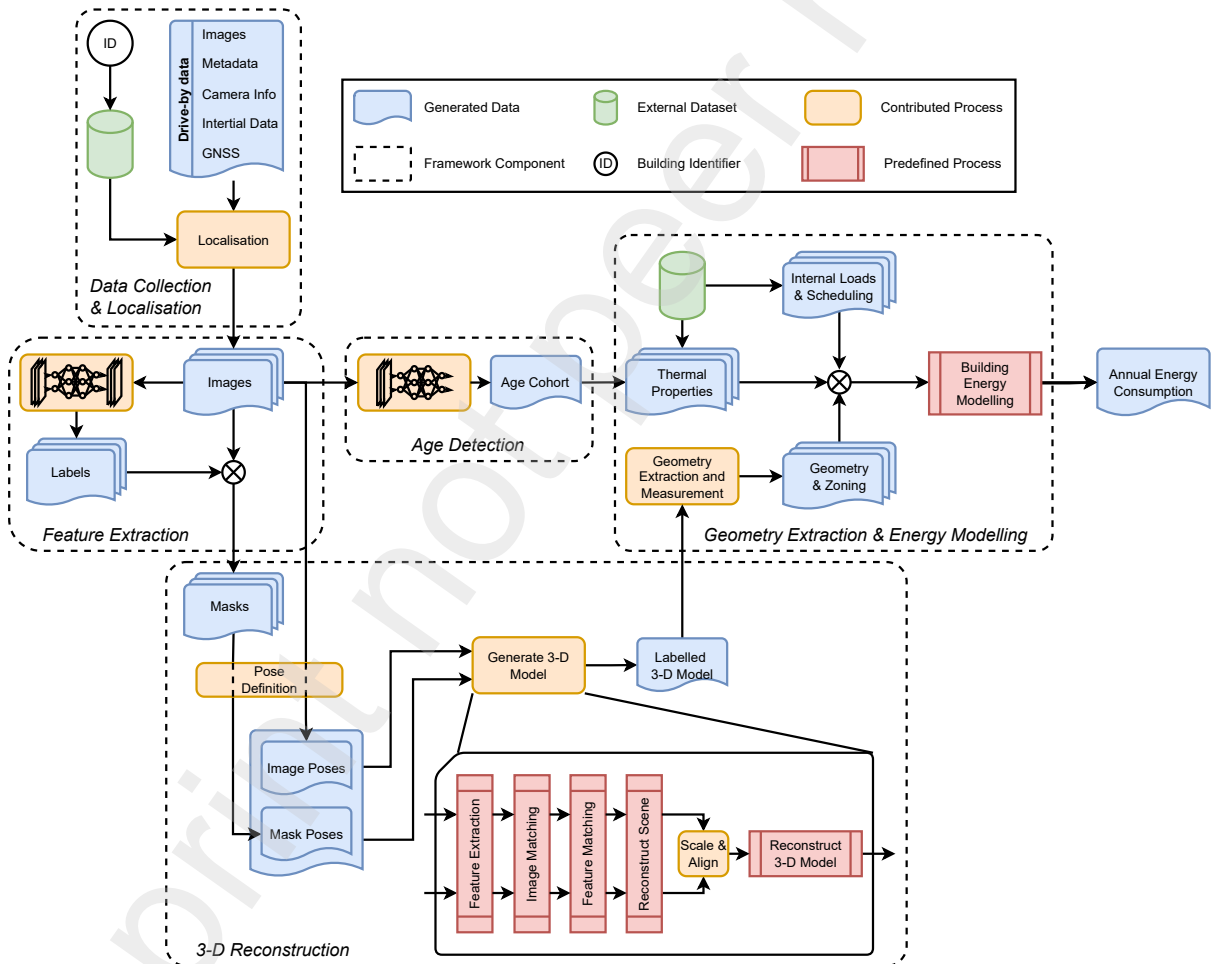


Figure 1: Overview of proposed energy consumption prediction framework, outlining individual components of the framework. Components of the framework are data collection and localisation (Section 2.1); age detection (Section 2.2); feature extraction (Section 2.3) and 3-D reconstruction (Section 2.4); and geometry extraction and energy modelling (Section 2.5). Aspects of the framework highlight where data is generated by the process, through either a process contributed by the authors or through predefined processes from publicly available software

79 and one upwards facing camera, each with a resolution of 5 megapixels (MP), capturing with
 80 a frequency of up to 30 frames per second (FPS). In practise, a trade-off is made between
 81 pixel resolution and capture frequency due to limitations of bandwidth in saving the images:
 82 uncompressed images comprise a huge amount of data so cannot be captured at very high
 83 frequencies. In this paper, higher resolution images are prioritised, with fewer high resolution
 84 images having been generally found to produce better quality results in the 3-D reconstruction
 85 component of the energy prediction framework; to this end, six 2048×2464 pixel images are
 86 captured at a rate of 10 FPS. Driving through a neighbourhood at approximately 4.5 m/s
 87 (≈ 16 km/h) around 12 images are captured per metre driven. At a distance of 10 m from the
 88 van, each pixel corresponds to approximately 2.5 cm^2 of, e.g., building facade.

89 A contemporary analogue to the image data produced is Google Street View (Anguelov et al.,
 90 2010), which has been used in both urban data projects (Li et al., 2017) and wider socioeconomic
 91 research (Gebru et al., 2017; Nguyen et al., 2020). Google Street View data is made available
 92 through a paid-for API. However, images are only available at a maximum pixel resolution of
 93 0.4 MP and are restricted in the available spatial resolution, with the API returning only the
 94 nearest image to a given location, which can limit details that can be extracted for a given
 95 property (Google, 2022).

96 2.1.2 Localisation

97 With the high volume of image data captured
 98 using a mobile sensing vehicle, a clear indexing
 99 scheme is required. Turning the images into a
 100 geospatial dataset requires robust localisation
 101 of the captured data, allowing images to be
 102 associated with a spatial pose which, in this
 103 paper, can be used to associate views with
 104 given houses. Onboard the sensing vehicle
 105 used, a georeferencing system² comprising an
 106 inertial measurement unit (IMU) and global
 107 navigation satellite system (GNSS) are used
 108 to monitor and map the location of the vehicle
 109 with an accuracy of up to 0.1 m. The IMU
 110 also provides the orientation of the vehicle, up
 111 to an accuracy of 0.1° . Localisation of the
 112 vehicle can be performed at a frequency of up
 113 to 100 Hz, equivalent to every hundredth of a
 114 second.

115 The onboard imaging system and IMU/GNSS share a time synchronisation system, which allows
 116 for reliable identification of the IMU/GNSS position for each image frame. With high frequency
 117 georeferencing of the van location, the position of each camera can be identified with linear

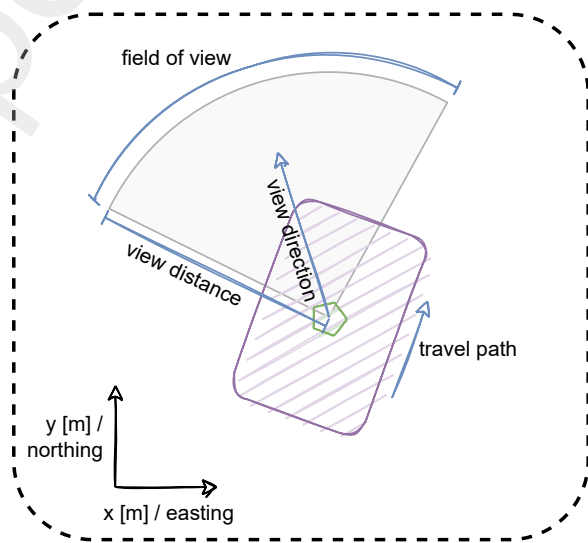


Figure 2: Sketch of a localised “view”, associated with an image, generated relative to the position and travel path of the sensing vehicle

²OxTS Survey+

118 interpolation of the vehicle’s position and orientation at capture time. At a driving speed of
119 4.5 m/s, there is an approximate accuracy assumption, for each camera frame’s position, of
120 0.25 m. The output position is recorded in World Geodetic System (WGS 84), i.e. longitude and
121 latitude. For use in the UK, the coordinates are reprojected into the Ordnance Survey National
122 Grid reference system (OSGB 1936), which gives the van’s position within the UK in metres,
123 allowing for direct measurement of 3-D models generated later in the framework.

124 Localising each frame is essential for extracting views of a given house. Individual “views” are
125 constructed to represent the perspective of a given image: the orientation and position of each
126 camera relative to IMU/GNSS unit, combined with the post-processed measurement of the
127 vehicle’s location and orientation will give the centre point of the camera, in OSGB 1936, and
128 the view direction of the camera. From this, a view is designed by creating a circular sector
129 from the absolute position of the camera, with some predefined view distance and field-of-view.
130 A sketch of this view is shown in Figure 2, highlighting one of the five radial views generated
131 from image data. The upward facing camera is disregarded in further processing.

132 As input to the framework, some building identifier is required to indicate the property of
133 focus. Two such identifiers commonly used for UK buildings are the unique property reference
134 number (UPRN) and Ordnance Survey topographic identifier (TOID). There are widely available
135 resources to link between these two references, as well as with other identifiers such as address
136 (Ordnance Survey, 2022a,c). With these identifiers, existing geospatial information of the
137 property such as its footprint can be extracted (Ordnance Survey, 2022b). Such information
138 can be used to associate localised images from drive-by capture with individual properties, by
139 finding intersections between the generated view and geospatial building information.

140 **2.1.3 Generating perpendicular views**

141 Due to the setup of the cameras, the images captured in the drive-by do not show perpendicular
142 “face-on” views of properties. For the age detection component of the framework, these views are
143 desirable as they provide a clear focal point for learning models. Such perspectives, however, can
144 be generated by reconstructing all images in a single frame as a panoramic image showing the
145 radial view of the sensing vehicle. Slicing this panorama can be used to create artificial views
146 to show new perspectives, including views perpendicular to the vehicle. However, these views
147 are not suitable for the reconstruction component, as they augment the images and remove
148 contextual information about the camera.

149 **2.2 Age Detection**

150 To estimate thermal properties of a building, statistical estimates of thermal transmittance, in
151 the form of u-values, are used due to the lack of comprehensive knowledge of thermal flux in
152 exterior features, such as walls, windows and roofs. The TABULA project provides estimates of
153 u-values based on age-based statistical archetypes for countries in the European Union, including
154 for Great Britain in the UK (BRE, 2014; Loga et al., 2016). Predicting the age-cohort will allow
155 the u-values from TABULA to be used to predict energy consumption.

156 **2.2.1 Model**

157 To classify the image data, a deep convolutional neural network (DCNN)-based model is used
 158 to estimate age cohort. Similar to the model presented in Zeppelzauer et al. (2018), the age
 159 detection model used in this paper relies on a patch-based classification and fusion approach,
 160 whereby the image is divided into subregions and each region is classified before an average
 161 pooling of the predicted age for each patch produces a single estimate.

162 The age detection model can be considered as three distinct elements: the feature selection,
 163 which extracts patches from the input image and stacks them into a single tensor; the backbone,
 164 a DCNN, that identifies and emphasises specific features within the patches; and a prediction
 165 layer, that distinguishes the features into classes, pooling classified patches to create a single
 166 unified prediction for the age cohort of the given image. Figure 3(a) shows an illustration of the
 167 model from input to prediction. The hyperparameters of the model are the number and size of
 168 patches, and the choice of backbone model. In the implementation of the framework, ResNet-18
 169 is used as a backbone, a common DCNN used in classification problems, with widely available
 170 implementations (He et al., 2016; PyTorch, 2022).

Table 1: Alignment of categories of age cohort used in Verisk UKBuildings, used to train the age detection model; TABULA archetypes used to estimate u-values, and aggregation reported in energy performance certificates (EPC)

UKBuildings	TABULA	EPC
Historic	pre-1919	before 1990
		1900-1929
Interwar	1919-1944	1930-1949
		1950-1966
Postwar	1945-1964	1967-1975
		1976-1982
Sixties Seventies	1965-1980	1983-1990
		1991-1995
Modern	1981-1987	1996-2002
		2003-2006
		2004-2009

171 **2.2.2 Training and Validation**

172 The age detection model was trained using a sample of 2463 images of houses in South Yorkshire,
173 UK, captured perpendicular to the mobile sensing vehicle, using the panorama slicing algorithm
174 described in Section 2.1.3. Building age cohorts were obtained using Verisk UKBuildings (Verisk,
175 2022), a geospatial dataset with a small number of attributes for residential properties in GB.
176 The aggregation of ages into cohorts is shown in Table 1, showing how the cohorts align with
177 the TABULA age categories, as well as the aggregation used in EPCs. While there is no perfect
178 alignment, the reliable availability of building ages in UKBuildings dictated its use as a label set
179 for training and validation.

180 The houses were randomly sampled from the capture data such that the number of houses in
181 each age cohort was approximately equal: 487 houses were ‘Historic’; ‘Interwar’ and ‘Postwar’
182 each comprised 496 houses; and 492 each of ‘Sixties Seventies’ and ‘Modern’. Houses for training
183 and validation were sampled from captured data in the South Yorkshire region, including
184 neighbourhoods in Sheffield and Barnsley, but excluding Doncaster as this is used in the case
185 study later.

186 The dataset was randomly subdivided into training, validation and testing sets, at a ratio of
187 80:10:10%, respectively. Training was performed initially for 50 epochs with early stopping using
188 validation loss at 32 epochs to prevent overfitting. Each epoch involved an evaluation of the
189 model, performed with a batch of $N = 8$ images before updating the parameters. Each image
190 was subdivided into $P = 32$ random 32×32 patches extracted from the middle 50% of the image,
191 and stacked and reshaped into an $(NP) \times 64 \times 64$ batch tensor before propagating through
192 the backbone and classification layers. The output predictions are reshaped and averaged to
193 produce a prediction for each class. The loss used cross entropy, and the Adam optimiser was
194 used with a learning rate of 0.003. The weights of the backbone were initialised with pretrained
195 weights for the classification of the ImageNet dataset to provide a well generalised starting point,
196 but were not fixed (Kornblith et al., 2019; PyTorch, 2022). All other weights in the model were
197 initialised randomly.

198 After 32 epochs of training, the validation accuracy was 71.2%, and the testing accuracy was
199 68.9%. The performance of the trained model is higher than that of the model proposed in
200 Zeppelzauer et al. (2018), but this could be accounted for by difference in the number of classes,
201 number of training, relative homogeneity of neighbourhoods in the training set and the different
202 approach to patching the images.

203 **2.3 Feature Extraction**

204 **2.3.1 Labelling Images**

205 Understanding images at a pixel level provides a means to identify features of a building facade. In
206 the framework, projecting labelled features into 3-D allows for the measurement and localisation
207 of individual components of the building. The action of labelling, or segmenting images, is to
208 assign each pixel to a set of semantic categories that inform the scene, essentially highlighting
209 what is in an image and where it is located. In the proposed framework, the features of interest

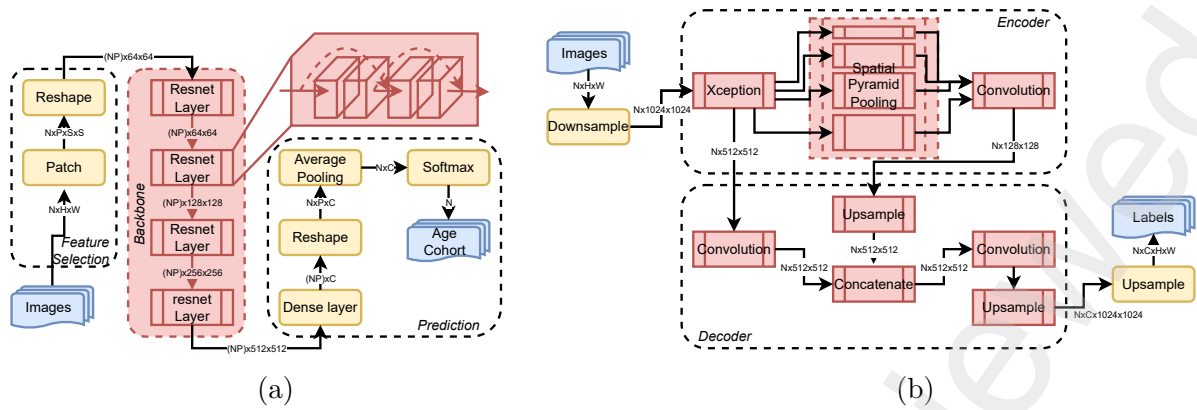


Figure 3: Illustrations of the two machine learning models used in the implementation of the proposed framework. (a) Age detection, using image patching to select features and pass through a backbone network, e.g. ResNet-18, and predict age cohorts based on average patch classification; (b) Feature extraction using semantic segmentation of images using the DeepLabv3+ architecture that encodes image features to different levels using a deep convolutional network, e.g. Xception, and spatial pyramid pooling, before decoding the features to pixel classes using convolution and concatenation of mixed-level features. Both illustrations indicate the dimensions of data throughout, where N is the number of $H \times W$ images processed, P is the number of $S \times S$ patches extracted, and C is the number of categories in the respective classification

210 are properties of the building facade and roof, namely windows, doors and chimneys, along
 211 with classification of the wall and roof. Semantic segmentation of building facades will return
 212 pixel-level labels of each category and treat anything else as “background”.

213 Manual segmentation of images is possible, but to do so on the scale required in the framework
 214 would be prohibitively time-consuming. To this end, machine learning approaches are utilised, as
 215 with the age detection component of the framework. Semantic segmentation of building facades
 216 is a well studied topic, with dedicated models designed around extracting building features (Dai
 217 et al., 2021; Ma et al., 2022). In this paper, we use DeepLabv3+, an out-of-the-box DCNN-based
 218 encoder-decoder model that is used for a wide range of semantic segmentation problems (Chen
 219 et al., 2018). The decision to use DeepLabv3+ was predominantly due to easily accessible
 220 implementations (PyTorch, 2022). The relative simplicity of the segmentation problem, in that
 221 facade features are typically simple rectangular shapes, lends to the idea that a generalised
 222 semantic segmentation model will perform well.

223 The DeepLabv3+ model is built around an encoder-decoder architecture. The encoder part
 224 generalises so-called “high-level” and “low-level” abstract features using a DCNN backbone, most
 225 commonly Xception (Chollet, 2017), and a series of algorithms called spatial pyramid pooling
 226 that learn to generate a low-dimensional representation of the models. These features are fed
 227 into the decoder part of the model, which learns a transformation to map these features to pixel
 228 level for classification. Training the model on a set of manually labelled building facades adapts
 229 the model to take in street-level images of houses and return a pixel map of semantic labels.
 230 Figure 3(b) shows an illustration of the basic architecture of the semantic segmentation model.

231 The model was trained using a set of 6000 directly captured images from the mobile sensing vehicle,
 232 which were manually annotated to highlight windows, walls, roofs, doors and chimneys. A small

233 proportion of the images used contained no labelled features, to provide better generalisation of
234 the trained model. The image was split 80:10:10% between training, validation and testing, and
235 was trained for 100 epochs. No early stopping was employed, but epochs were capped at 100 due
236 to limits in computational resources. The test accuracy of the model, i.e. the average percentage
237 of pixels correctly classified, was 93.6%. A more discerning metric, the mean intersection over
238 the union (IOU) across all labels, was recorded as 78.9%. The IOU quantifies the degree of
239 overlap between predicted regions and true segmented regions, and is widely used in classification
240 problems (Dai et al., 2021). The results for the trained model are in line with state-of-the-art
241 semantic segmentation work, e.g. (Dai et al., 2021; Ma et al., 2022). Figure 4 shows the results
242 of the trained model alongside the ground truth for an example image in the test set.

243 2.3.2 Masking Images

244 The trained segmentation model is used to automatically create label maps for facades to be
245 used for projection and measurement in the framework. An additional benefit of these label
246 maps is that they can be used to mask the original images to remove background features, which
247 is beneficial during 3-D reconstruction, as the final model will only contain features belonging
248 to a building, without additional objects like cars or other urban furniture such as trees and
249 lampposts. Reconstructing only the building in the images reduces the amount of post-processing
250 required to extract geometries from the 3-D model. An example of a masked image is shown in
251 Figure 4.

252 2.4 3-D Reconstruction

253 Once a set of views of a building has been labelled and masked, the 3-D reconstruction component
254 of the framework is performed. Using the known localised views, as described in Section 2.1.2,
255 along with intrinsic properties of the cameras, such as focal distance and field-of-view, poses can

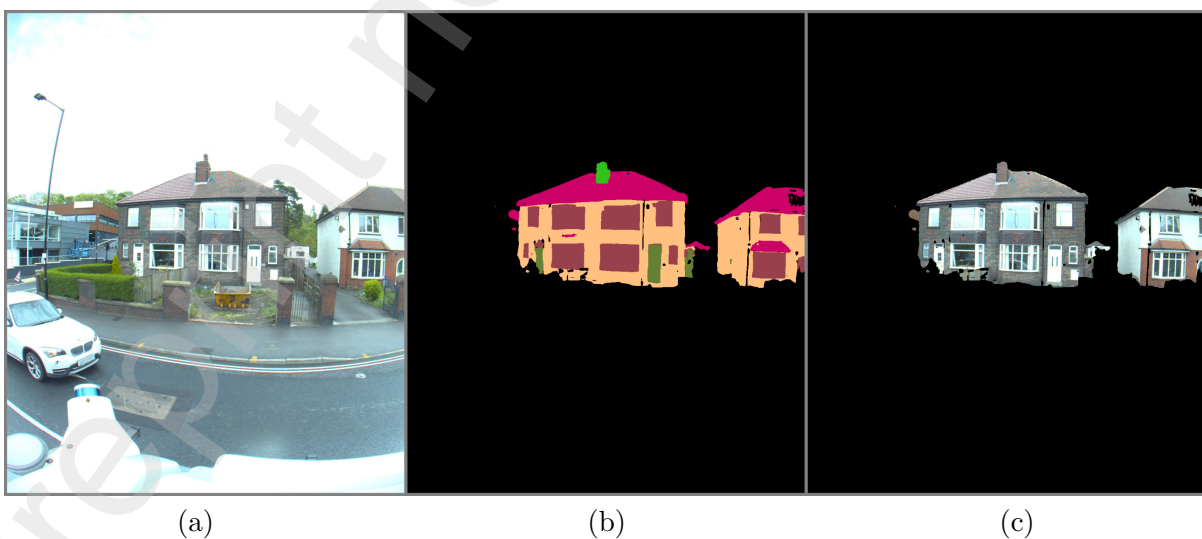


Figure 4: Demonstration of labelling and masking of a building facade in a drive-by captured image. (a) The raw image; (b) A label mask resulting from evaluation of the trained semantic segmentation model; (c) The masked image obtained by removing “background” features

256 be reconstructed to build 3-D models with real-world coordinates. Mapping image and label
257 data onto these 3-D models also allows for specific facade features to be measured in real-world
258 coordinate space for use in energy modelling.

259 **2.4.1 Defining Poses**

260 Generating poses from the data is required to codify the spatial information about each image.
261 In the framework, for a given building, there is a geospatially located polygon representing its
262 global position within GB in metres. Also localised are the position and ‘view’ of associated
263 images that contain the given building, identified during the localisation step. To generate
264 intermediary data used for 3-D reconstruction, positions are centred relative to the centroid of
265 the polygon, by simply translating the global position of the images to be positioned relative to
266 the polygon, allowing for easier measurement of the output reconstruction.

267 Additionally, intrinsic camera properties such as the focal length and lens distortion, as well as
268 the orientation in 3-D space of the camera, are attributed to each image pose. This process is
269 repeated for both the original images and the masked images to create two sets of reconstruction
270 data.

271 **2.4.2 Generating 3-D Building Model**

272 Once a set of poses is generated from the labelled and masked images, the 3-D reconstruction
273 component of the framework can be used to build geometric models of buildings. Using a
274 combination of structure-from-motion and multi-view stereoscopy, the multiple perspectives of
275 the building can be used to localise features and create a surface model in 3-D space (Griwodz
276 et al., 2021). The first step in the generation of a 3-D model is the detection and extraction of
277 the so-called “features” in each image. These features are identified using the scale-invariant
278 feature transform (SIFT), a widely used algorithm that detects abstract descriptive properties in
279 an image, based on various properties such as sudden changes in colour or shape. The invariance
280 of the features allows them to be paired together regardless of any perceptive transformation
281 or distortion they are affected by, e.g. rotation, translation or shearing (Lowe, 2004). SIFT
282 features have been used across computer vision applications, including object recognition in
283 video tracking and image stitching, as well as 3-D reconstruction (Peng et al., 2009). The ability
284 of SIFT to provide a generalised representation of features in images allows for the pairing and
285 matching of objects to create correspondences in the building facade and wider urban furniture
286 in the localised images.

287 The list of features extracted from each image is used to pair images based on their relative
288 poses. Typically, in a structure-from-motion pipeline, this process relies on finding common SIFT
289 descriptors between images and assigning pairings based on matches. However, the localisation
290 gives known poses which allows for simpler assumptions to be made on the association between
291 images. The assumption made here is that each feature can only have one corresponding match,
292 which reduces the computational requirements that impact processing time, but limits the
293 effectiveness, especially on repetitive structures. Despite this, the process remains fairly robust.
294 To improve matching, this process is repeated twice, once for the original images and once for the

Table 2: Information and parameters required for energy modelling

Geometry and Zoning	Coordinates of footprint
	Building height
	Window coordinate locations
	Number of storeys
Thermal Properties	Thermal transmittance of wall, floor, roof and window
	Solar heat gain coefficient of windows
	Outdoor air filtration
Meteorological conditions	Weather data
Internal loads and scheduling	Intensities of occupancy, lighting and equipment
	Schedules of occupancy, lighting and equipment
	Building service systems

295 masked images. The original images contain many more features within an image that can be
 296 used to infer context in the scene, including objects such as trees and cars that increase the total
 297 number of matches to improve the quality of image and feature referencing. However, since the
 298 intended output of the reconstruction is just the building, applying the process to the masked
 299 images generates a secondary representation of the scene. An alignment of the mask-based scene
 300 with the full feature reconstruction, in conjunction with the predefined poses acts as a corrective
 301 transformation that minimises potential errors and inaccuracies in, e.g. the location of the van,
 302 and the lesser contextual information in the masked images.

303 With the aligned and matched features, a reconstruction of the 3-D model can be performed
 304 by creating a surface mesh by connecting features. The generated mesh is automatically post-
 305 processed to remove artefacts and reduce regions with a large number of nodes. Following this, a
 306 texture map is created, which essentially projects the images onto the reconstructed 3-D model.
 307 Applying the texturing process with both the masked images and the label maps provides two
 308 representations of the house: one with photographic detail, and the other with a semantic label
 309 localised in 3-D space. In the former case, this can be used for visualisation, and checking the
 310 quality of a reconstruction, while the projected labels allow for the extraction and measurement
 311 of geometry of the building.

312 2.5 Geometry Extraction and Energy Modelling

313 To scale up the (partial) automation of building energy models requires condensing and formatting
 314 of building properties needed to simulate consumption over a defined period of time. In this
 315 section, the individual considerations of the building model are described and processes to extract
 316 them from drive-by imaging, or otherwise, are detailed.

2.5.1 Geometry Extraction and Measurement

Building a process to automate the extraction of geometry used to build energy models requires some preparation of the 3-D models generated in the pipeline.

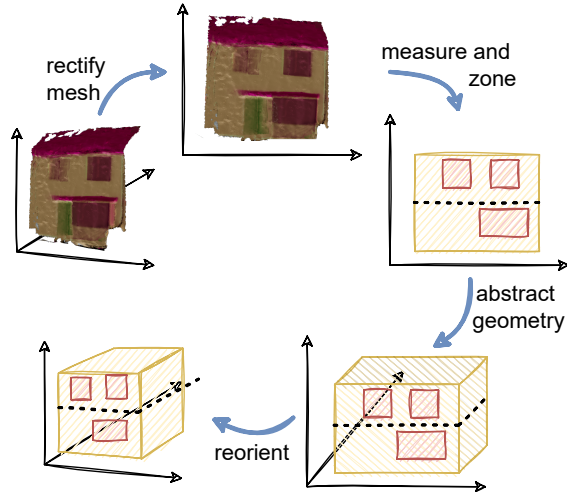


Figure 5: Illustration of pipeline for building geometric representation of building from a labelled reconstructed mesh, showing the process of rectifying a mesh to orient with the yz -plane; measuring the facade with bounding boxes and zoning based on storeys; abstracting the facade to a 3-D geometry; and transforming the building to its original orientation.

The orientation of the mesh is transformed to best align with the unit axes, such that the front facade aligns with the yz -plane. Bounding boxes are fitted to the features on the facade, such as the wall and each disconnected window, based on the projected labels. These bounding boxes are used to generate the measured geometry used in the building energy models. The geometry is separated into zones based on the number of storeys, which is inferred from the orientation of windows. The generated geometric representation is retransformed to the mesh's original orientation. Figure 5 illustrates the geometry extraction and measurement component of the framework.

2.5.2 Building Energy Modelling

The culmination of the building information extracted from drive-by images is to construct a model from which energy consumption, and therefore performance, can be effectively simulated. In this work, EnergyPlus, an industry-standard whole building energy simulation program, is used to estimate energy consumption given the inputs from the framework. These inputs are defined in an intermediary data format, generated from the information described in Table 2.

This approach allows for the modelling of each building independently, providing energy consumption information at a high level of granularity. The core aspects of the building that are identified as important for the simulation are the geometry and zoning; the thermal properties; and the internal loads and schedules. In the case of the former two, data obtained from drive-by captured data is used.

Due to the automated nature of the framework, the mesh generated during 3-D reconstruction may contain artefacts or low quality regions. A preprocessing step to remove low-quality features is performed by identifying natural clusters in the 3-D model: the DBSCAN algorithm identifies distinct spatial regions based on areas with dense detail (Ester et al., 1996); by virtue of the semantic segmentation and masking, largest distinct cluster is considered to represent the building, and the rest is discarded.

The orientation of the mesh is transformed to best align with the unit axes, such that the front facade aligns with the yz -plane. Bounding boxes are fitted to the features on the facade, such as the wall and each disconnected window, based on the projected labels. These bounding boxes are used to generate the measured geometry used in the building energy models. The geometry is separated into zones based on the number of storeys, which is inferred from the orientation of windows. The generated geometric representation is retransformed to the mesh's original orientation. Figure 5 illustrates the geometry extraction and measurement component of the framework.

356 **Geometry and Zoning** Creating the physical representation of a building is, in essence, the
357 extrusion of its footprint using the calculated geometry. Thermal zones can be inferred from
358 the number of storeys, obtained by counting windows. The zones, representing a storey, are
359 assumed to be of equal height, for simplicity in the generalisation of zones. The windows are
360 represented as coordinated quadrilaterals on the facade, extracted from the minimum bounding
361 rectangle for each window on the 3-D reconstructed mesh. For non-visible faces of the building,
362 the window-to-wall ratio is instead encoded: equal to that of the measured facade on the opposite
363 face; and a low number, e.g. 10%, for side faces. Symmetry of the window-to-wall ratio and
364 a low non-zero number for side faces were used in the absence of directly observed data, as
365 reasonably considered assumptions of the average construction of a residential building. A visual
366 representation of the process is shown in Figure 5.

367 **Thermal Properties** To infer the thermal properties of the building, a set of age-based
368 typologies for GB, developed by BRE, were used to infer the u-values of different features.
369 The TABULA age-cohorts, which were used in the age detection component of the proposed
370 framework, contain statistical assumptions of u-values for different properties, including walls
371 and windows. The estimates given by the age detection model were used to generate estimates
372 for the thermal transmittance and solar heat gain coefficient to be input in the energy model.

373 **Internal Loads and Scheduling** Due to the lack of observable information on the internal
374 properties of each building, a uniform assumption was made for all simulations. To keep these
375 assumptions as close to those used in EPCs as possible, the heating schedules were sourced from
376 SAP 2012 guidelines: 9 hours on week days and 16 hours on weekends (BRE, 2012); schedules
377 were also modelled to approximately represent reported diurnal patterns in energy usage (Few
378 et al., 2022). Lighting and electrical equipment, and occupancy scheduling, were referenced from
379 literature and industry guidelines (Baden-Powell et al., 2011; BRE, 2012).

380 **2.5.3 Simulating Energy Consumption**

381 In addition to the systems defined, associated weather data was obtained from International
382 Weather for Energy Calculations (IWE, 2001). Based on the extracted geometry and generated
383 model, energy consumption based on space heating, lighting and equipment is simulated using
384 EnergyPlus giving an estimation of annual usage in kWh/m².

385 **3 Related Work**

386 Data driven solutions to categorising and quantifying the built environment, particularly energy
387 consumption, are numerous and long-standing (Fouquier et al., 2013). Much of the research
388 has focussed on understanding material stocks and predicting energy performance at large scales
389 (Arbabi et al., 2021; Wei et al., 2018). To this end, building automated energy models from
390 building data has been researched, for example in Chen et al. (2017) and Wate and Coors (2015).
391 In Chen et al. (2017), the authors develop scenarios for retrofit at city-scale, using building data.
392 3-D data models have also been used to simulate energy usage (Wate and Coors, 2015).

393 Detection of facade features using machine learning and computer vision has become a popular
394 topic in the last few years (Chen et al., 2017). Tailor-made facade segmentation solutions, such
395 as in Liu et al. (2017) and Dai et al. (2021) report high accuracy but are limited in that they
396 are applied predominantly to rectified images, i.e. those that have had lens distortion features
397 removed, similar to the format used in the age detection component of the framework. Due to
398 a lack of code availability, and specific requirements for the format of images, neither solution
399 was used in this work. Other features that have been identified from street-level images include
400 building age (Zeppelzauer et al., 2018), and heating energy demand (Despotovic et al., 2019).

401 Another component of the proposed framework relies on the projection of features to 3-D for
402 the extraction of geometry. Identifying properties of buildings using existing sources such
403 as Google Street View (Anguelov et al., 2010) has been applied to improving understanding
404 of the urban environment (Hara et al., 2013; Campbell et al., 2019). Feature detection and
405 mapping from Google Street View has been used to estimate building heights and improve facade
406 understanding (Yuan and Cheriyyadat, 2016). One of the main limitations with Google Street
407 View data, however, is the spatial and temporal resolution at which it is available, meaning it
408 can be difficult to reconstruct high quality 3-D geometries.

409 Aerial remote sensing is has also been used in urban quantification: datasets such as UKBuildings
410 utilise LiDAR to estimate building heights (Verisk, 2022); Bayomi et al. (2021) use thermography
411 to calibrate building envelopes; and remote sensing has been used to develop material stock
412 models (Lanau et al., 2019). Stock models of buildings, such as in London, have been developed
413 to build representations of cities, utilising data sources including aerial LiDAR (Steadman et al.,
414 2020).

415 **4 Case Study: Doncaster, UK**

416 The framework as outlined was applied to houses in a residential neighbourhood in Doncaster,
417 UK. Data from 53 council-owned social houses were selected for the case study, based on available
418 information including addresses. Publicly available energy performance certificates (EPC) were
419 obtained for each of the properties and filtered for the most up-to-date version. The framework
420 was applied to each house and the estimated features of the building, including height and energy
421 consumption, are compared against existing data sets.

422 **4.1 Existing Data**

423 The three main datasets used to compare and validate the outputs of the framework are OS
424 MasterMap, Verisk UKBuildings and EPC reports. In the case of the former two datasets,
425 properties such as the building footprint, building height and building age are available, with
426 varying degrees of quality and confidence. In EPC reports, predictions of the energy consumption,
427 alongside features used in the inputs of the SAP modelling are provided, including dwelling
428 features such as heating systems, room space, and assumed insulation.

429 While polygons from OS MasterMap were used to localise the data, all other available properties
430 are kept separate. In practice, one might combine these sources of available data to gain greater

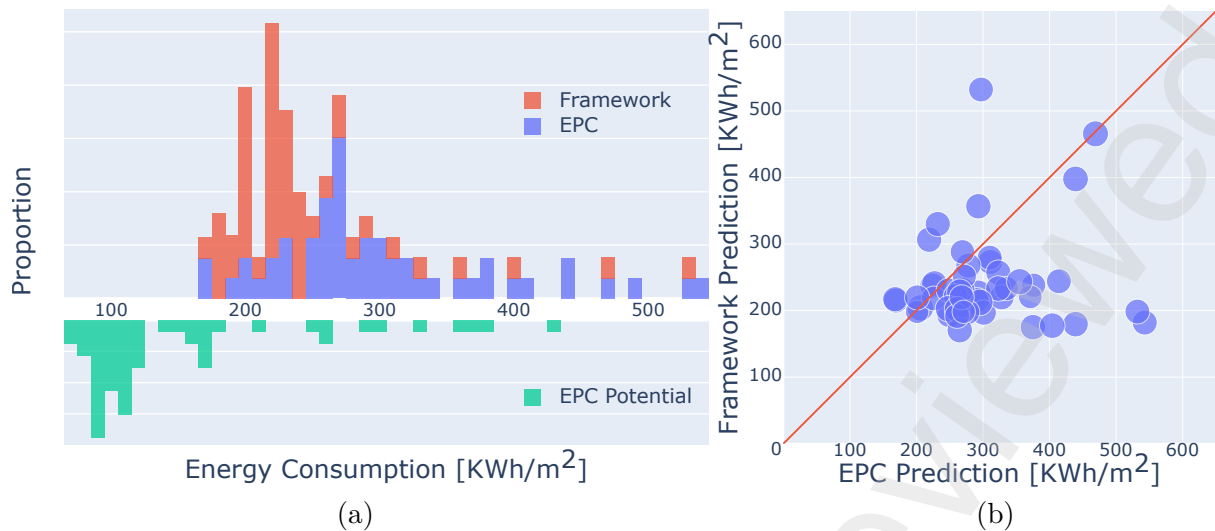


Figure 6: Comparison of predicted energy consumption by the MARVel-based framework and EPCs. (a) the distribution of predicted energy consumption from SAP-based EPCs and from the proposed framework. The lower plot shows the distribution of “potential” energy consumption based on recommendations provided in the EPC reports. (b) the energy consumption predictions from each source for each property in the case study.

431 representation of a building, but this is beyond the scope of the paper.

432 4.2 Validation

433 The predicted annual energy consumption in EPC reports is provided in kWh/m^2 , which is
 434 used to generate ratings. Using the energy modelling component of the framework, annual
 435 energy consumption, in kWh/m^2 , can be simulated. In these models, simulation takes less than
 436 a minute on a mid-tier laptop.

437 Figure 6 shows the predicted energy consumption by each method. In Figure 6(a), the distribution
 438 of energy consumption predictions over the sample houses is shown, highlighting the similarities
 439 in the overall values predicted across the houses estimated. The bottom plot indicates potential
 440 energy consumption provided in the EPC based on retrofit recommendations, to give a visual
 441 reference for the potential range of energy consumption values. In Figure 6(b), the energy
 442 consumption by each method for each building is shown, for direct comparison of estimation.
 443 From Figure 6(a), there is demonstrable agreement in the estimations for the majority of
 444 houses, with the framework-based estimation tending towards slightly lower estimations of
 445 energy consumption, on average. This is corroborated in the scatter plot, with most properties
 446 predicted as equal or slightly lower energy consumption with the framework. In some cases,
 447 however, there is a larger difference between the predictions, either with the EPC reporting
 448 higher values for energy consumption, or vice versa. In such cases, these can largely be attributed
 449 to the assumptions made both in the generation of the energy model using the drive-by data,
 450 which considers all internal systems and scheduling to be uniform across all houses; and in
 451 the assumptions made by the EPC provider. For example, for one property, the EPC-based
 452 energy consumption has been predicted significantly higher than the framework-based approach.
 453 Looking at the report features, the EPC highlights issues with poor efficiency from windows,

454 walls and the water heating system, the latter of which forms the most significant aspect of
455 the recommendation to reduce energy consumption. On the other hand, the framework-based
456 approach has assumed an average internal condition due to the lack of other data made available
457 to it, and assumed thermal properties based on statistical archetypes of the building, characterised
458 by its predicted age.

459 To validate the geometry extracted from the
460 3-D reconstructed models, the difference in
461 calculated height with data available from a
462 3rd party dataset for each house is shown in
463 Figure 7. The relative height-to-the-eaves, re-
464 ported as “RelH2” in OS MasterMap Build-
465 ing Height Attribute (OS, 2022), is used as
466 a benchmark, as the building heights are not
467 reported in the EPC data. Figure 7 shows a
468 general agreement between the two estimates
469 of facade height, with a mean difference of
470 0.245 m.

471 The case study conducted has demonstrated
472 the feasibility of street-level drive by capture,
473 in that it has the capability to provide estim-
474 ates that are largely in agreement with those
475 reported by EPCs, without the need for entry
476 to the property for direct inspection. Factor-
477 ing in the degree of assumptions also made in
478 EPCs on aspects such as insulation thickness,
479 the drive-by approach has viability, even with
480 the approximations of interior state required.

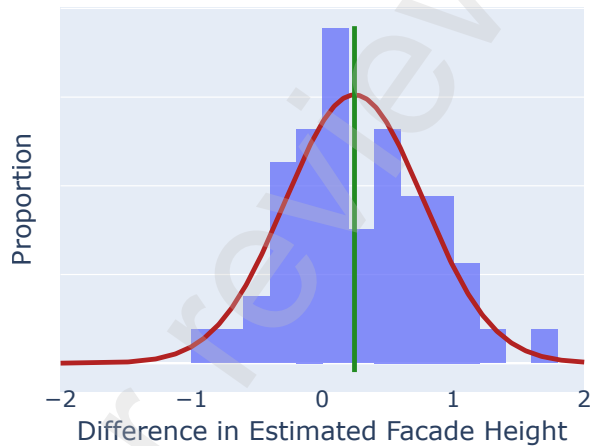


Figure 7: Histogram showing the proportion of differences between the facade predicted in the framework and the value reported in OS MasterMap for the buildings in the case study. When the difference is greater than zero, this indicates the framework height is greater than the OS value. The mean difference is 0.245 m highlighted by the vertical line. A Gaussian curve fitted to the mean and standard deviation is overlaid.

481 5 Limitations, Modularity and Extensibility

482 The framework proposed and outlined in this paper is designed to perform large-scale generation
483 of digital representations of buildings, with the end-goal of the framework to measure energy
484 consumption. There are a number of limitations with each component of the framework, as in
485 most data-driven approaches, including the requirements for large scale deployment and the
486 limited availability of high quality ground truth data for validation.

487 These are discussed in turn in the following subsections. However, the modularity of the
488 framework allows for clear adaptation and extension to both overcome some of the identified
489 limitations and facilitates application to wider research problems in the urban environment.

490 5.1 Framework Design

491 **Data Collection and Localisation** A source of uncertainty is the localisation of the sensing
492 vehicle, which can be located to within 0.1 m, but in practice may be less accurate. Furthermore,
493 the difference in projection of location used between different sources of data, and the conversion
494 between, is also a potential source of error or inaccuracy. While the GNSS/IMU unit in the van
495 measures the global position in longitude and latitude, with units in degrees, building identifiers,
496 such as OS TOIDs, as well as the measurement of structures, are made in metres. The conversion
497 from longitude and latitude to easting and northing used by Ordnance Survey has an inaccuracy
498 up to 1 m, and this conversion is only available in third-party GIS softwares.

499 The availability of other sources of data is limited, even for nationwide datasets. For example,
500 UKBuildings only has building age for 71% of the houses in the neighbourhood from which the
501 case study buildings were taken, and there is no clear way to validate the accuracy of these
502 values.

503 It may also be possible to expand on the information provided from drive-by data collection
504 by including additional data modalities, such as thermal and LiDAR. Greater insight and
505 reliability of the structure of the buildings that these modalities might add will lead to more
506 confident estimations of the energy consumption, as well as providing additional inferences, such
507 as insulation thickness or fault detection by assessing thermal properties and more accurate
508 geometric structure using the LiDAR point clouds.

509 Independent of uncertainty, another limitation is the sheer volume of data obtained by drive-by
510 capture. In contrast to parameter based datasets, where each house is represented as a set
511 of variables, typically text- or number-based, the drive-by process of capturing images and
512 geolocation information creates a huge amount of data that needs to be processed and stored.
513 While this paper outlines the framework as a proof-of-concept, the quantity of computational
514 storage and other resources required need to be considered before deploying this type of image-
515 based solution at a city, regional or national scale. In terms of raw image data, 1000 images takes
516 up approximately 1 GB space. At a capture rate of 7.5 FPS, the total storage for a 75-minute
517 drive is around 150 GB.

518 **Feature Extraction** Other sources of uncertainty in the framework include the trained
519 machine learning models used to identify age and component features. In the latter case, most
520 feature pixels are classified with an accuracy well over 95%; an exception to this are roofs, which
521 are classified with accuracy 81% – the source of uncertainty here may be in the varied inclusion
522 of eaves and gutters in the training data. For the age detection part of the framework, the
523 model is only accurate approximately 70% of the time. While this is in line with state-of-the-art
524 models for detection, it will likely have an impact on the results.

525 Despite the relatively high levels of accuracy in the models used in this work, they are not
526 necessarily the best performing of all available models. As mentioned previously, bespoke facade
527 detection methods might result in higher quality label data, and as research progresses, the
528 state-of-the-art will improve. As the framework is designed to be modular in terms of defined

529 inputs and outputs, replacing components with new methods should be simple. However, with
530 the addition of multiple modalities, such as thermal data, joint representation learning could
531 further extend the capabilities of feature extraction: in the identification and localisation of
532 facade features; or providing additional insight into the thermal properties of the building
533 (Theodosiou et al., 2021).

534 **Geometry and Zoning** When creating the geometric representation of the building relies
535 on a number of simplified assumptions. Thermal zoning is assumed by storey and, due to the
536 lack of consistent reconstruction of roofs, all buildings are assumed to have a flat roof even if
537 the contrary is true. As shown in the results, neither of these assumptions have drastic effects
538 on estimated energy consumption, compared to the values reported in EPCs. Incorporating
539 additional knowledge, for example structural archetypes, to better generalise the assumptions
540 may yield more representative estimations of energy consumption.

541 The 3-D reconstruction component of the framework is based on multiview stereoscopy, but
542 using different approaches to photogrammetry may yield more accurate measurements. LiDAR,
543 for example, measures the 3-D scene directly, and aerial point clouds have been used for urban
544 quantification (Steadman et al., 2020). More contemporary methods in 3-D scene representation
545 include neural radiance fields (NeRF), which utilise generative machine learning models to
546 generalise views and poses (Mildenhall et al., 2020). NeRFs have, for example, been used to
547 represent cities at different spatial scales (Xiangli et al., 2022).

548 **Thermal Properties** Detection of materials from visual images alone is difficult due to issues
549 such as paint and texture, so the framework uses an age detection model and statistical values
550 for the thermal transmittance for use in the energy modelling component. However, combined
551 information from data sources and drive-by data might be able to infer u-values with greater
552 resolution. For example, whether a facade wall is solid brick, or a cavity wall, filled or unfilled,
553 can be used to estimate u-values with greater degree of accuracy using industry documentation,
554 such as reduced data SAP (rdSAP) (BRE, 2012). How to infer such properties with drive-by
555 data is an open question that requires further research.

556 **Internal Loads and Scheduling** As highlighted in the evaluation of the case study, it is
557 difficult to quantify the internal loads and scheduling used for energy modelling assumptions
558 from drive-by capture. There are many studies which cover modelling occupant behaviour that
559 can be used to provide “best guesses” (Carlucci et al., 2020; Yang et al., 2021). Focusing on
560 this is beyond the scope of the paper, but introducing seasonality in the internal loads and
561 scheduling, based on trends identified from qualitative research, or from statistical models based
562 on smart meter data could introduce a more realistic model of occupant behaviour (Few et al.,
563 2022).

564 **5.2 Validation and Verification**

565 On validation of the approach, there is limited data, especially on the interior of the property,
566 that is consistently reliable and available. EPCs, UKBuildings and OS Mastermap all use similar

567 data driven approaches, or rely on flawed assessment as in the case of EPCs, and due to this
568 there is no robust information for verifying data (Hardy and Glew, 2019).

569 Validation of the geometry measured with the framework is limited due to the availability of
570 consistent and reliable data. In the case study, building height calculated in the framework was
571 compared with OS data. However, for every house in the case study, the data is considered
572 unverified by Ordnance Survey. In fact, this is the case for almost all houses in the area local to
573 the case study. For example, of over 382,000 houses analysed in South Yorkshire, UK, 99.65%
574 were reported as unverified in the OS MasterMap Building Height Attribute dataset (OS, 2022).

575 **5.3 Scalability**

576 Consumer advice sites estimate that conducting an EPC survey can take between 45 and 60
577 minutes per house (ReallyMoving.com, 2022). With drive-by capture, up to 30 houses can
578 be imaged per minute, based on the assumptions made in this paper. While the current
579 implementation of the framework was built to model and simulate energy consumption for a
580 single house at a time, it is possible to expand the scale of the framework by incorporating parallel
581 computing, to reconstruct and simulate energy for multiple houses simultaneously. Similarly,
582 additional efficiencies might be found in preprocessing of the image data. Captured data contains
583 multiple houses and could be batch processed, including with labelling and masking.

584 **5.4 Beyond Energy Consumption**

585 While the focus in this paper has been on using drive-by capture data to model energy con-
586 sumption, there are many other uses for the data and the processes discussed. Understanding
587 material inventories and dimensions can help build up a picture of material stock (Lanau et al.,
588 2019). An extension of the proposed framework, with statistical assumptions based on factors
589 such as age and build factor, would be to build a database to help quantify material stock for a
590 specific region (Arbabi et al., 2021).

591 Similar to understanding material intensity in a region, knowledge of the dimensions and
592 retrofit needs of a building can yield solutions to facilitate efficient manufacturing for retrofit
593 interventions, e.g. for mass-produced panelised systems (Orlowski, 2020). Quantifying housing
594 stocks at a neighbourhood, town or city level would allow for efficient manufacturing, resulting
595 in an economy of scale benefit (Andronie et al., 2021).

596 **6 Conclusions**

597 This paper has outlined a multi-aspect, modular framework spanning from the capture and
598 localisation of drive-by image data to reliable, scalable energy consumption prediction of
599 individual residential buildings. The resulting predictions produce similar estimates for energy
600 consumption as to EPCs. Each component of the framework is discussed and critically appraised
601 both in terms of its individual performance and its contribution to the wider output of the
602 framework. Limitations and extensions to the framework are discussed and more general aspects
603 of the process are explored.

604 The current implementation of the feature extraction and measurement aspect of the pipeline
605 is designed to measure geometry and simulate energy consumption for a single property. This
606 design allows for parallel execution when scaling up to neighbourhood- or city-level. While
607 reliant on the same input data set of drive-by capture, each reconstruction is independent,
608 computational requirements notwithstanding. As opposed to reconstructing and measuring
609 whole streets, the highly parallelisable approach in our methodology is much better designed for
610 future scaling.

611 In practice, a hybrid method is likely needed, one that takes into account data from multiple
612 sources and scales, with a view to confidence in data, its quality and any uncertainty in the
613 process. The modular aspect of the framework presented in this framework offers some facilitation
614 for improvements as methods improve, and for extensions into new avenues.

615 Authorship Contribution Statement

616 **Wil O C Ward:** Conceptualization, Methodology, Software, Validation, Formal analysis,
617 Investigation, Writing – Original Draft, Funding acquisition. **Xinyi Li:** Conceptualization,
618 Methodology, Software. **Yushu Sun:** Data Curation, Software. **Menglin Dai:** Software. **Hadi**
619 **Arbabi:** Conceptualization, Writing – Review & Editing, Project administration, Funding
620 acquisition. **Danielle Densley Tingley:** Writing – Review & Editing, Supervision, Funding
621 acquisition. **Martin Mayfield:** Supervision, Funding acquisition.

622 Declaration of Competing Interest

623 The authors declare that they have no known competing financial interests or personal relation-
624 ships that could have appeared to influence the work reported in this paper.

625 Acknowledgements

626 This work was supported by EPSRC Active Building Centre [EP/V012053/1]. WOCW, XL and
627 HA were additionally supported by Towards Turing 2.0 under EPSRC [EP/W037211/1] and
628 The Alan Turing Institute. Neither EPSRC nor The Alan Turing Institute had any involvement
629 in study design; execution; or in the writing of this article.

630 References

631 International weather for energy calculations. *American Society of Heating, Refrigerating and*
632 *Air-Conditioning Engineers, ASHRAE*, 2001.

633 M. Andronie, G. Lăzăroiu, M. Iatagan, I. Hurloiu, and I. Dijmărescu. Sustainable cyber-physical
634 production systems in big data-driven smart urban economy: a systematic literature review.
635 *Sustainability*, 13(2):751, 2021.

636 D. Anguelov, C. Dulong, D. Filip, C. Frueh, S. Lafon, R. Lyon, A. Ogale, L. Vincent, and

- 637 J. Weaver. Google street view: Capturing the world at street level. *Computer*, 43:32–38, 6
638 2010. doi: 10.1109/MC.2010.170.
- 639 H. Arbabi, M. Lanau, X. Li, G. Meyers, M. Dai, M. Mayfield, and D. D. Tingley. A scalable data
640 collection, characterization, and accounting framework for urban material stocks. *Journal of*
641 *Industrial Ecology*, 2021. doi: 10.1111/JIEC.13198.
- 642 C. Baden-Powell, J. Hetreed, and A. Ross. *Architect’s pocket book*. Oxford: Architectural P.,
643 Oxford, 4th ed. edition, 2011. ISBN 9780080969596.
- 644 N. Bayomi, S. Nagpal, T. Rakha, and J. E. Fernandez. Building envelope modeling calibration
645 using aerial thermography. *Energy and Buildings*, 233:110648, 2021.
- 646 BRE. The Government’s standard assessment procedure for energy rating of dwellings, 2012.
647 URL <https://www.gov.uk/guidance/standard-assessment-procedure>.
- 648 BRE. Building typology brochure England, 2014.
- 649 A. Campbell, A. Both, and Q. C. Sun. Detecting and mapping traffic signs from google street
650 view images using deep learning and gis. *Computers, Environment and Urban Systems*, 77:
651 101350, 2019.
- 652 S. Carlucci, M. De Simone, S. K. Firth, M. B. Kjærgaard, R. Markovic, M. S. Rahaman, M. K.
653 Annaqeeb, S. Biandrate, A. Das, J. W. Dziedzic, et al. Modeling occupant behavior in
654 buildings. *Building and Environment*, 174:106768, 2020.
- 655 L. C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous
656 separable convolution for semantic image segmentation. *Lecture Notes in Computer Science*
657 *(including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*,
658 11211 LNCS:833–851, 9 2018. doi: 10.1007/978-3-030-01234-2_49. URL [https://link.
659 springer.com/chapter/10.1007/978-3-030-01234-2_49](https://link.springer.com/chapter/10.1007/978-3-030-01234-2_49).
- 660 Y. Chen, T. Hong, and M. A. Piette. Automatic generation and simulation of urban building
661 energy models based on city datasets for city-scale building retrofit analysis. *Applied Energy*,
662 205:323–335, 11 2017. doi: 10.1016/j.apenergy.2017.07.128. URL [https://linkinghub.
663 elsevier.com/retrieve/pii/S0306261917310024](https://linkinghub.elsevier.com/retrieve/pii/S0306261917310024).
- 664 F. Chollet. Xception: Deep learning with depthwise separable convolutions. *2017 IEEE*
665 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017-January:1800–1807, 7
666 2017. doi: 10.1109/CVPR.2017.195.
- 667 M. Dai, W. O. Ward, G. Meyers, D. D. Tingley, and M. Mayfield. Residential building facade
668 segmentation in the urban environment. *Building and Environment*, 199:107921, 7 2021. doi:
669 10.1016/J.BUILDENV.2021.107921.
- 670 M. Despotovic, D. Koch, S. Leiber, M. Doeller, M. Sakeena, and M. Zeppelzauer. Prediction
671 and analysis of heating energy demand for detached houses by computer vision. *Energy and*
672 *Buildings*, 193:29–35, 2019.

- 673 M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A density-based algorithm for discovering
674 clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996.
- 675 J. Few, M. Pullinger, E. McKenna, S. Elam, E. Webborn, and T. Oreszczyn. Smart energy
676 research lab: Energy use in GB domestic buildings 2021. *Smart Energy Research Lab (SERL)*
677 *Statistical Reports*, 1, 2022.
- 678 A. Foucquier, S. Robert, F. Suard, L. Stéphan, and A. Jay. State of the art in building modelling
679 and energy performances prediction: A review. *Renewable and Sustainable Energy Reviews*,
680 23:272–288, 2013.
- 681 T. Gebru, J. Krause, Y. Wang, D. Chen, J. Deng, E. L. Aiden, and L. Fei-Fei. Using deep
682 learning and Google Street View to estimate the demographic makeup of neighborhoods across
683 the united states. *Proceedings of the National Academy of Sciences*, 114(50):13108–13113,
684 2017.
- 685 Google. Street View Static API, 2022. URL [https://developers.google.com/maps/
686 documentation/streetview](https://developers.google.com/maps/documentation/streetview). [Online; accessed: 11-July-2022].
- 687 C. Griwodz, S. Gasparini, L. Calvet, P. Gurdjos, F. Castan, B. Maujean, G. D. Lillo, and
688 Y. Lanthony. Alicevision meshroom: An open-source 3d reconstruction pipeline. ACM Press,
689 2021. doi: 10.1145/3458305.3478443.
- 690 J. Gwak, J. Jung, R. D. Oh, M. Park, M. A. K. Rakhimov, and J. Ahn. A review of intelligent
691 self-driving vehicle software research. *KSII Transactions on Internet and Information Systems*
692 (*TIIS*), 13:5299–5320, 11 2019. doi: 10.3837/TIIS.2019.11.002.
- 693 K. Hara, V. Le, and J. Froehlich. Combining crowdsourcing and google street view to identify
694 street-level accessibility problems. pages 631–640, 2013.
- 695 A. Hardy and D. Glew. An analysis of errors in the energy performance certificate database.
696 *Energy policy*, 129:1168–1178, 2019.
- 697 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings*
698 *of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 699 IPCC. Climate change 2022: Mitigation of climate change, 2022.
- 700 S. Kornblith, J. Shlens, and Q. V. Le. Do better imagenet models transfer better? In *Proceedings*
701 *of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2661–2671,
702 2019.
- 703 M. Lanau, G. Liu, U. Kral, D. Wiedenhofer, E. Keijzer, C. Yu, and C. Ehlert. Taking stock of
704 built environment stock studies: Progress and prospects. *Environmental Science & Technology*,
705 53(15):8499–8515, 2019. doi: 10.1021/acs.est.8b06652.
- 706 X. Li, C. Zhang, and W. Li. Building block level urban land-use information retrieval based on
707 Google Street View images. *GIScience & Remote Sensing*, 54(6):819–835, 2017.
- 708 H. Liu, J. Zhang, J. Zhu, and S. C. Hoi. Deepfacade: A deep learning approach to facade

709 parsing. *IJCAI International Joint Conference on Artificial Intelligence*, 0:2301–2307, 2017.
710 doi: 10.24963/IJCAI.2017/320.

711 T. Loga, B. Stein, and N. Diefenbach. Tabula building typologies in 20 european countries—mak-
712 ing energy-related features of residential building stocks comparable. *Energy and Buildings*,
713 132:4–12, 2016.

714 D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of*
715 *Computer Vision* 2004 60:2, 60:91–110, 11 2004. doi: 10.1023/B:VISI.0000029664.99615.94.
716 URL <https://link.springer.com/article/10.1023/B:VISI.0000029664.99615.94>.

717 W. Ma, S. Xu, W. Ma, and H. Zha. Multiview feature aggregation for facade parsing. *IEEE*
718 *Geoscience and Remote Sensing Letters*, 19, 2022. doi: 10.1109/LGRS.2020.3035721.

719 B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. NeRF:
720 Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.

721 Q. C. Nguyen, Y. Huang, A. Kumar, H. Duan, J. M. Keralis, P. Dwivedi, H.-W. Meng, K. D.
722 Brunisholz, J. Jay, M. Javanmardi, et al. Using 164 million Google Street View images to
723 derive built environment predictors of COVID-19 cases. *International journal of environmental*
724 *research and public health*, 17(17):6359, 2020.

725 Ordnance Survey. AddressBase, 2022a. URL [https://www.ordnancesurvey.co.uk/](https://www.ordnancesurvey.co.uk/business-government/products/addressbase)
726 [business-government/products/addressbase](https://www.ordnancesurvey.co.uk/business-government/products/addressbase). [Online; accessed: 12-July-2022].

727 Ordnance Survey. MasterMap topography layer, 2022b. URL [https://www.ordnancesurvey.co.](https://www.ordnancesurvey.co.uk/business-government/products/mastermap-topography)
728 [uk/business-government/products/mastermap-topography](https://www.ordnancesurvey.co.uk/business-government/products/mastermap-topography). [Online; accessed: 12-July-
729 2022].

730 Ordnance Survey. Open Linked Identifiers, 2022c. URL [https://www.ordnancesurvey.co.uk/](https://www.ordnancesurvey.co.uk/business-government/products/open-linked-identifiers)
731 [business-government/products/open-linked-identifiers](https://www.ordnancesurvey.co.uk/business-government/products/open-linked-identifiers). [Online; accessed: 12-July-
732 2022].

733 K. Orłowski. Automated manufacturing for timber-based panelised wall systems. *Automation in*
734 *construction*, 109:102988, 2020.

735 OS. Mastermap building height attribute, 2022. URL [https://www.ordnancesurvey.co.uk/](https://www.ordnancesurvey.co.uk/business-government/products/mastermap-building)
736 [business-government/products/mastermap-building](https://www.ordnancesurvey.co.uk/business-government/products/mastermap-building).

737 K. Peng, X. Chen, D. Zhou, and Y. Liu. 3d reconstruction based on sift and harris feature
738 points. *2009 IEEE International Conference on Robotics and Biomimetics, ROBIO 2009*,
739 pages 960–964, 2009. doi: 10.1109/ROBIO.2009.5420735.

740 PyTorch. Models and pre-trained weights, 2022. URL [https://pytorch.org/vision/stable/](https://pytorch.org/vision/stable/models.html)
741 [models.html](https://pytorch.org/vision/stable/models.html). [Online; accessed: 13-July-2022].

742 ReallyMoving.com. How to get an EPC, 2022. URL [https:](https://www.reallymoving.com/energy-performance-certificates/guides/how-to-get-an-energy-performance-certificate-epc)
743 [//www.reallymoving.com/energy-performance-certificates/guides/](https://www.reallymoving.com/energy-performance-certificates/guides/how-to-get-an-energy-performance-certificate-epc)
744 [how-to-get-an-energy-performance-certificate-epc](https://www.reallymoving.com/energy-performance-certificates/guides/how-to-get-an-energy-performance-certificate-epc). [Online; accessed: 23-August-
745 2022].

- 746 L. D. Shorrocks, J. Henderson, and J. I. Utley. Reducing carbon emissions from the uk housing
747 stock. 2005. URL www.bre.co.uk.
- 748 P. Steadman, S. Evans, R. Liddiard, D. Godoy-Shimizu, P. Ruyssevelt, and D. Humphrey.
749 Building stock energy modelling in the uk: the 3dstock method and the london building
750 stock model. *Buildings and Cities*, 1:100–119, 5 2020. doi: 10.5334/BC.52. URL <http://journal-buildingscities.org/articles/10.5334/bc.52/>.
751
- 752 T. Theodosiou, K. Tsikaloudaki, K. Kontoleon, and C. Giarma. Assessing the accuracy of
753 predictive thermal bridge heat flow methodologies. *Renewable and Sustainable Energy Reviews*,
754 136:110437, 2021.
- 755 UKGBC. The Retrofit Playbook. Technical report, 2021.
- 756 UNEP. 2020 global status report for buildings and construction: Towards a zero-emission,
757 efficient and resilient buildings and construction sector, 2020.
- 758 Verisk. UKBuildings, 2022. URL <https://digimap.edina.ac.uk/verisk>. [Online; accessed:
759 12-July-2022].
- 760 P. Wate and V. Coors. 3d data models for urban energy simulation. *Energy Procedia*, 78:
761 3372–3377, 2015.
- 762 Y. Wei, X. Zhang, Y. Shi, L. Xia, S. Pan, J. Wu, M. Han, and X. Zhao. A review of data-driven
763 approaches for prediction and classification of building energy consumption. *Renewable and*
764 *Sustainable Energy Reviews*, 82:1027–1047, 2018.
- 765 Y. Xiangli, L. Xu, X. Pan, N. Zhao, A. Rao, C. Theobalt, B. Dai, and D. Lin. Bungeenerf:
766 Progressive neural radiance field for extreme multi-scale scene rendering. In *The European*
767 *Conference on Computer Vision (ECCV)*, 2022.
- 768 B. Yang, F. Haghghat, B. C. Fung, and K. Panchabikesan. Season-based occupancy prediction
769 in residential buildings using machine learning models. *e-Prime-Advances in Electrical*
770 *Engineering, Electronics and Energy*, 1:100003, 2021.
- 771 J. Yuan and A. M. Cheriyyadat. Combining maps and street level images for building height and
772 facade estimation. pages 1–8, 2016.
- 773 M. Zeppelzauer, M. Despotovic, M. Sakeena, D. Koch, and M. Döllner. Automatic prediction of
774 building age from photographs. In *Proceedings of the 2018 ACM on International Conference*
775 *on Multimedia Retrieval*, pages 126–134, 2018.