# Generating a nationwide residential building types dataset using machine learning

Kristina Dabrock [a,b,*] , Jens Ulken [a], Noah Pflugradt [a], Jann Michael Weinand [a] , Detlef Stolten [a,b]

[a] *Forschungszentrum Jülich GmbH, Institute of Climate and Energy Systems, Jülich Systems Analysis, Jülich, 52425, Germany*
[b] *RWTH Aachen University, Chair for Fuel Cells, Faculty of Mechanical Engineering, Aachen, 52062, Germany*

## ARTICLE INFO

## ABSTRACT

The lack of high-resolution building data is an obstacle to the development of detailed, spatially explicit recommendations for decarbonization measures. In an effort to fill this gap, this study outlines the creation of a building level dataset based on standardized building archetypes for all German residential buildings. A machine learning approach using XGBoost is used to train models to predict the size class and construction year of individual buildings. Refurbishment states are assigned based on federal state level statistics. Based on these characteristics, TABULA building archetypes are assigned. The training data generation is primarily based on the grid dataset of the German census. The data is enriched with morphological features of buildings and neighborhoods, as well as socio-economic characteristics. The machine learning models perform with accuracies of 97.4 % and 73.9 %, respectively, on a test set at the individual building level. The distribution of size classes and construction years in the resulting dataset shows a high degree of agreement with official statistics at the federal state level, but also a tendency to overrepresent majority classes. This study proves that the chosen methodology is suitable for generating a complete nationwide dataset. By providing spatially resolved, individual building data that can serve as a proxy for the energetic properties of buildings, the resulting dataset can facilitate building-related energy transition analyses.

## 1. Introduction

Buildings contribute significantly to energy consumption and greenhouse gas emissions in both the EU [1] and Germany [2]. In order to devise targeted measures with the goal of reaching greenhouse gas-neutrality, detailed knowledge of the building stock is required. Statistical data about the building stock and its energy-related statistics is available for European countries on a national level, e.g., through the Building Stock Observatory [3] or JRC-IDEES [4], and on a regional level, for example for Germany [5]. An increasing amount of open data initiatives are publishing data on the building level, but detailed, structured, complete, and accessible building level information, even for geometric data [6], let alone going beyond, is still scarce. This data, however, is necessary to carry out spatially explicit analyses that allow assessments of, e.g., heat demand at the local level.

A common approach for handling the scarcity of building level data in bottom-up analyses is enriching existing building data with data from building typology archetypes, as in Schwanebeck et al. [7] and Yang et al. [8]. These typologies, ranging from the national (e.g., TABULA [9]) to regional (e.g., building typology for the German federal state Schleswig–Holstein [10]) levels, classify buildings into categories based on a set of characteristics. The TABULA typology [9] defines building categories through a code composed of letters and numbers indicating location, size class, construction year, and refurbishment state of a building (see Fig. 1). Apart from Single Family Houses (SFHs), the typology differentiates between Multi-Family Houses (MFHs), Terraced Houses (THs) and Apartment Blocks (ABs). While the typology dates to 2012, it remains to date the most comprehensive typology for categorizing the German residential building stock known to the authors and is therefore an important basis for the large-scale analysis of heat demand and decarbonization strategies. A challenge lies in missing categories for modern buildings, constructed after the publication of the typology. However, it has even been updated after its initial publication and includes a construction year class beginning in 2016 [11]. Furthermore,

---

**Region**:                    **Construction year**
National                          **period**:
(no regional specification)    1860-1918

D E . N . S F H . 0 2 . G e n . R e E x . 0 0 2

**Country**:          **Type**:            **Refurbishment state**:
Germany    Single Family              Usual
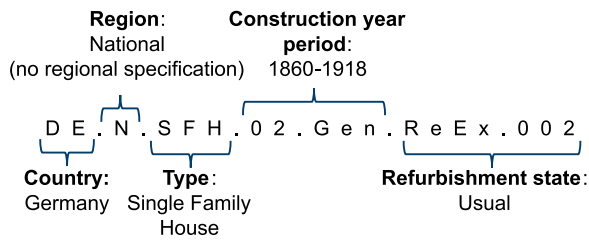                    House

**Fig. 1.** TABULA typology code explanation and example.

buildings built after 2012 make up only 5 % of the residential building stock [12], and it can be argued that these are the least relevant buildings in terms of energy refurbishment measures due to the already higher energy efficiency standards. And while discrepancies between U-values in TABULA and energy performance certificates raise concerns regarding the reliability of the TABULA typology [13], classifying buildings according to their archetypes is still useful, since it currently is the best data available for Germany and can furthermore easily be updated, once better data becomes available. Additional data such as, for example, heat transmission values, areas of components, and specific heat demand, is available for each of the archetypes that represent a typical building of that category. Although the data provided by the TABULA typology naturally does not exactly match each building falling into a certain category, it provides a valuable approximation of energy-related characteristics. These characteristics are essential for performing large-scale building level energy simulations, which in turn are required for the assessment of decarbonization potentials and scenarios. A prerequisite for harnessing the data provided in the TABULA typology in spatially explicit analyses is the assignment of TABULA types (see Fig. 1) to individual buildings. This requires data on construction year, size class, and the refurbishment state of a building. Leveraging the assigned TABULA type, simulation models for individual buildings such as ETHOS.HiSim [14] can be configured by accessing U-values and component areas provided by the TABULA typology. These models can then be employed for analyzing, for example, the large-scale heat demand of residential buildings in Germany [15].

This study presents a methodology for assigning building archetypes from the TABULA typology [9] to residential buildings. For the first time, such a methodology is developed to analyze the building stock of the entire country of Germany. A machine learning-based approach was employed for determining the archetype-defining characteristics. As a result, the dataset includes not only the TABULA type but also the construction year, size class, and refurbishment state of each building, giving users the flexibility to use these attributes in a variety of possible analyses. Finally, the generated dataset was validated against official statistics and a detailed assessment of its quality was carried out. The complete and validated TABULA dataset for Germany provides a valuable basis for spatially explicit research on building energy demand.

## 2. Related work

Several previous studies have developed methods for assigning size classes to buildings. Yang et al. [8], for instance, use building level data on the number of shared walls, the number of registered addresses, building footprint areas, gross floor areas, and the number of stories to define a rule-based mapping of size classes (SFH, mid-TH, end-TH, apartment building or MFH) to Dutch buildings. A similar rule-based mapping, based on the ground floor area and number of floors, was also implemented by Schwanebeck et al. [7] and Blanco et al. [16], with 250 m² being a common threshold for differentiating between SFHs and MFHs. Wurm et al. [17] and Droin et al. [18] both use Random Forest classifiers trained on data from homogeneous census grid cells for predicting one of three or four size class types, respectively. The former article presents a case study for the city of Münster, Germany, and

reaches an overall accuracy of 96 %. The latter study applies the model to all of the federal states in Germany, with an overall accuracy of more than 95 %.

In addition to the size classes, the year of construction also carries valuable information about a building. The construction year is not only necessary to derive the TABULA type but can serve as a proxy for heat demand-related characteristics such as insulation [19], level heights, fabric, and glazing and is also relevant beyond energy research, e.g., for natural hazard risk assessments [20]. As the individual buildings' construction year is not available for many countries, this data must be gathered, e.g., from local building cadasters, commercial datasets, remote sensing data, or occupant surveys [7]. Approaches for detecting a building's construction year are still under development [21].

Addressing the issue of manual data collection, Alexander et al. [22] present automated approaches for deriving construction year information from map data using cluster analysis based on building shape and context information. For Germany, statistical data is available down to the municipality and hectare levels from the census of 2011 [23]. The use of this data has been shown to lead to acceptable results for heat demand modeling at the city level by Zirak et al. [24], who applied and evaluated their methodology for two small German towns. Thus, a simple approach for assigning construction year periods to individual buildings consists of assigning the predominant construction year period within a census grid cell to all buildings within that grid cell, as demonstrated in Wurm et al. [17]. Although this is accurate for homogenous grid cells, it leads to errors in non-homogenous ones.

Machine learning can be applied to take individual building and context characteristics into account. Algorithms such as Convolutional Neural Networks (Zeppelzauer et al. [25] and Li et al. [26]) and Random forest (Biljecki and Sindram [27], Rosser et al. [20], Garbasevschi et al. [19], and Blanco et al. [16]) are common choices, but the input data and features vary. Zeppelzauer et al. [25] reach accuracies of 55.1 % by using a Convolutional Neural Network trained on building photographs for predicting the decade of construction. Meanwhile, Li et al. [26] report a mean absolute error of 11 years when predicting construction years combining a Convolutional Neural Network and Support Vector Regression and training on image data from Google Street View. Biljecki and Sindram [27] train a Random Forest model with nine features derived from a 1-D building dataset of the city of Rotterdam, Netherlands. The study reaches a mean absolute error of between 4.9 and 19.4 years, depending on the available attributes. Rosser et al. [20] train a Random Forest model for predicting five construction year periods based on 15 building morphology and neighborhood characteristics for Notthingham, United Kingdom, and reach accuracies of 77 %. Similarly, Garbasevschi et al. [19] predict ten construction year periods using a Random Forest model based on more than 50 building, street and block metrics as features from open spatial data for eight cities in the German federal state North Rhine–Westphalia with accuracies of up to 80 % [19]. The authors find that the spatial distribution of training and test data influences the accuracy, and that training and test buildings should ideally be in close spatial proximity, increasing the accuracy up to 96 % if training data from the same city is used [19].

The third aspect required to determine the TABULA types is data on the refurbishment status of a building. A lack of refurbishment state data at the building level has been highlighted by Zirak et al. [24] and Yang et al. [8]. Even when data is gathered at the building level, an assignment to individual buildings is problematic due to privacy issues [24]. Both Zirak et al. [24] and Yang et al. [8] resort to randomly distributing refurbishment states to individual buildings based on higher-level statistics. Wurm et al. [17] circumvent the issue of missing refurbishment data by defining scenarios for the potential refurbishment states of all buildings.

A straightforward mapping of building attributes to TABULA types based on building footprint polygons, construction years, and building heights is presented by Yang et al. [8]. In this study, the authors outline a framework for modeling residential space heating energy demand using

GIS data in combination with TABULA archetypes for the city of Leiden, Netherlands. However, due to limited data availability in many other regions, this approach is not easily transferable. Wurm et al. [17] also model heat demand for a city in Germany by assigning TABULA types; however, their assignment of construction years and size classes is simplified and refurbishment states are not assigned. Schwanebeck et al. [7] calculate the heat demand of the northern German region of Schleswig–Holstein by combining geographic 3-D building data with census datasets. However, although they use 3-D building data as an input, their output is heat demand on the hectare level, thus not providing TABULA types for individual buildings.

In summary, existing studies on the assignment of size classes, construction years, refurbishment states, and TABULA types have either a limited spatial scope, lack a methodology that is suitable for generating nationwide datasets, or have not been validated against official statistics. In this study, this gap is filled by providing TABULA types for all residential buildings in Germany using an individual building-based methodology. Contrary to past studies, all characteristics required for deriving the TABULA type are assigned and the methodology is applied to the entire country of Germany.

## 3. Methodology and data

The methodological structure of this paper is illustrated in Fig. 2. Section 3.1 describes the input data. The underlying building data used is geographic and was extracted from OpenStreetMap and from official governmental 3-D building datasets for Germany (see Section 3.1.3). Combined with census (see Section 3.1.1) and socio-economic data (see Section 3.1.4), this constitutes the training data used for training XGBoost machine learning models (see Section 3.2). These models were then applied to predict the construction year and size class of all residential buildings in Germany. Furthermore, the refurbishment state of buildings was assigned probabilistically based on the statistics provided by the German Environment Agency [28]. Based on these three characteristics, TABULA types were assigned to each individual building (see Section 3.3), providing access to a wide range of additional information available for each type through the TABULA typology, such as insulation levels, window areas, and heat demand.

### 3.1. Input data

The main input data sources used in this study are census data, refurbishment state data, individual building data and socio-economic data. The census data is relevant for creating the training dataset by extracting buildings that can be labeled with a construction year and size class. Census data, individual building data, and socio-economic datasets serve as additional sources for deriving features that the machine learning model can then be trained on. The probabilistic refurbishment state assignment is based upon federal state level statistics about refurbishment states.

### 3.1.1. Census data

The census data is used for training label generation in order to classify the construction year and size classes, as well as feature generation. The German census [29] constitutes a statistical examination intended for the aggregation of population and housing data. It contains data such as construction year classes, size classes, inhabitants, and the number of households. Each data entry represents the composition of a surveyed attribute of a geographically-referenced grid cell, sized 100 m x 100 m [30]. This approach is applied to maintain anonymity and circumvent potential deductions regarding individual identities. Furthermore, in order to inhibit potential deanonymization, data reported within a cell can be deliberately distorted in cases where the inputs for a particular grid cell are insufficient.
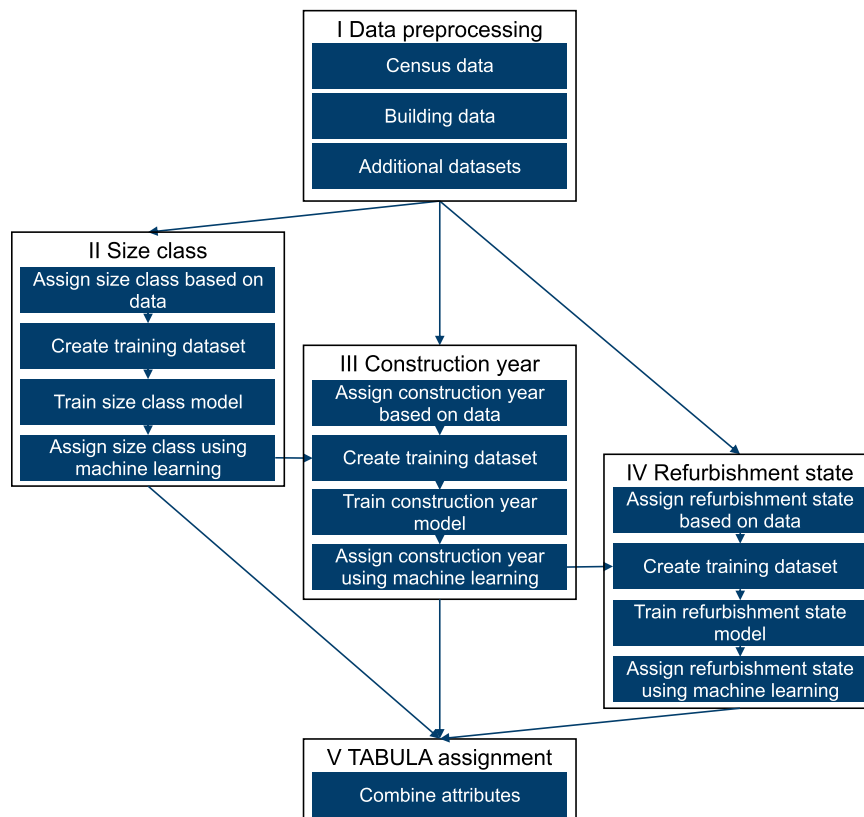


**Fig. 2.** Flow chart illustrating the data processing and assignment workflow.

### 3.1.2. Refurbishment state data

As the census data does not contain information on the refurbishment state of buildings, another source is required. No reliable sources with building level data for the assignments of refurbishment states could be identified. Therefore, the probabilistic assignment of refurbishment states is based on statistical data available from the German Environment Agency [28]. They report the share of buildings in the three categories "not refurbished", "partially refurbished", and "fully refurbished", as well as the share of new buildings, at the federal state level.

### 3.1.3. Building data

To achieve the goal of assigning TABULA types to individual buildings, features for the machine learning model must be derived from individual building data.

The basic building data, including building footprints and height, was extracted from official government 3-D building datasets at Level of Detail 2 (LoD2), i.e., including standardized roof shapes, for the German federal states where it was available at the time of carrying out this study (Bavaria, Berlin, Brandenburg, Hamburg, Hesse, Lower Saxony, North Rhine–Westphalia, Saxony, Saxony–Anhalt, and Thuringia), and for all other federal states from OpenStreetMap, data was retrieved from the Geofabrik download server [31]. This process and the resulting data are described in detail in Dabrock et al. [6].

Buildings were classified by their usage type as *residential, non-residential*, or *mixed*. Training data was generated by implementing a direct mapping from the function attribute and OpenStreetMap tags in the LoD2 data and OpenStreetMap data, respectively. Only the type information in Saxony-Anhalt was disregarded as it included implausible information. As discussed by Bandam et al. [32], the building use of the majority of buildings in Germany is not specified in OpenStreetMap. At the time of writing this study, 68.28 % of German buildings are tagged as "building='yes'", instead of providing more detailed information [33], leading to data gaps which require data imputation steps. Based on the training data, a machine learning model was trained (for details refer to Supplementary Material S1) and then applied to assign types to all previously unclassified buildings. In the following steps of this study, only the buildings classified as *residential* or *mixed* are included.

Additionally, buildings were enriched with data about the roof, including the following roof characteristics: type, height, orientation, tilt, and area. From the 3-D building data, the *roofType* attribute of a building as defined in the ALKIS Objektartenkatalog [34] and extended by SIG 3-D for CityGML [35], which was extracted and mapped to a selection of the roof shapes defined in OpenStreetMap [36]. Height was calculated by subtracting the height of the lowest point of a roof section from the height of its highest point and taking the average over all roof parts. The azimuth, tilt, and areas of all roofs were calculated using functions from the *polygon3dmodule* from Biljecki et al. [37]. For the buildings originating from OpenStreetMap, the roof characteristics of shape, height, tilt, and orientation were read from tags "*roof:shape*," "*roof:height*," "*roof:angle*," "*roof:direction*," respectively. However, only 7.7 % of German buildings have information for "*roof:shape*", and for the other attributes, the share remains below 1 % [38]. Roof shape values were cleaned by mapping to more common, similar categories and removing those that cannot be mapped. Valid roof shape categories and mapping tables are included in Supplementary Material S2. Height and tilt data was cleaned by selecting only the first part of the value if a space is present and removing the strings "m," "°,""deg," and "~" in order to remove units and other modifiers. Values not matching the floating point format after cleaning were ignored. Textual representation of roof orientation was mapped to degrees with "N" for north being set to 0 [°].

The footprint area, i.e., the contact area of a building with the ground, can be directly calculated from the basic building data by calculating the area of the footprint polygon. However, in the context of heat demand modeling, more relevant measures can be calculated by taking the number of floors or, as a proxy, the building height into

consideration. The details of the calculation are included in Supplementary Material S3.

### 3.1.4. Socio-economic data

The training dataset was further enriched with socio-economic data. Data pertaining to socio-economic factors at the municipal and county levels were accumulated based on the assumption that a positive socio-economic environment may serve as an indicator for relatively recent construction years and higher refurbishment frequencies for buildings. Studies by the German Economic Institute and analyses of the market-leading German real estate portal Immoscout24, show a negative correlation between price discounts and energy efficiency classes of buildings, while modern, energy-efficient buildings are stable in price [39,40].

Assuming that higher acquisition costs presuppose a high economic status, in the case of MFHs, these acquisition costs lead to higher rents, which results in the gentrification of neighborhoods and therefore socio-economic data may allow for inferences about housing stock quality. As Edlund et al. [41, p. 23] state, "Gentrification is about price growth and changes to the housing stock, not population growth." It is further presumed that higher income and lower unemployment rates could be predictive of SFHs, given the prerequisite of a high net worth or income for procurement and maintenance. The collected spatial socio-economic data is mimicked after the attributes within the GRW-Indicator (Indicator for the Improvement of the Regional Infrastructure) [42]. In the German context, the GRW-Indicator is a tool utilized by the Federal Ministry for the Economy and Climate Protection to discern socio-economically-underprivileged regions and correspondingly allocate financial support [42, pp. 1–6]. The calculation of this indicator involves the use of regional gross annual wages per employee, regional unemployment rates, projections of regional employment in correlation with overall German development, and the Infrastructure Indicator 2012 [43]. The Infrastructure Indicator 2012 provides a picture of Germany's physical capital-, human capital-, and household-oriented infrastructure. The physical capital-oriented infrastructure is composed of three sub-indicators, which display the accessibility of the three nearest national or foreign conurbations, the equipment with high-level transport infrastructure, and the level of high-performance broadband infrastructure. The human capital-oriented infrastructure is described by apprenticeship capacities, the number of employees in knowledge-intensive and business-oriented services, the number of employees in technical professions, and in knowledge transfer institutions. Lastly, the regional population potential represents the household-oriented infrastructure [43]. Beyond socio-economic data, this study also employs RegioStaR17 [44] data. RegioStaR17, a regional statistical spatial typology, stratifies geographical regions into 17 distinct types such as metropolises, urban areas, small town regions, and rural areas, amongst others [44].

### 3.2. Machine learning approach for size class and construction year assignment

The input data previously presented was used to create a labelled dataset for training and testing, and to derive features for training machine learning models. Then, the actual model training process, including hyperparameter optimization, was carried out.

### 3.2.1. Labelling data for model training

As in Garbasevschi et al. [19] and Droin et al. [18], target labels for construction year and size classes were generated by assigning the values of census grid cells containing only one attribute expression to all buildings located in the respective grid cell, e.g., cells that contain only SFHs. Furthermore, a specific parameter, indicative of the level of anonymization, is attributed to each grid cell. Consequently, only grid cells exhibiting the lowest degree of anonymization were retained. This serves to ensure maximum data quality for the generated training

dataset. The grid cells in the census dataset have a size of 100 m x 100m. The assigned attribute expressions to a grid cell describe the building type topology of the buildings located within it. Therefore, attributes are not available on a single building level. By labeling only buildings located in grid cells with one attribute expression for all buildings in the grid cell, combined with the lowest degree of anonymization, the possibility of a correct assignment of the building type on a single building level is maximized. All buildings that could be labelled using this procedure were retained for the model training process.

The TABULA framework categorizes buildings into four size classes, compared to ten found in the census data. Therefore, following the methodological approach described by Loga et al. [45], a mapping (see Table 1) was applied to assign each building in the census data one of the four size classes defined by TABULA. This was performed before extracting the grid cell data.

The labelled dataset pertaining to building size class classification contains 10.8 million observations (see Table 2), which accounts for 55.7 % of the German residential building stock in 2021 [46, p. 19], each characterized by 41 features. The SFH size class dominates the label distribution, accounting for approximately 86 % of instances, as shown in Fig. 3. As can be seen, ABs are severely underrepresented. The class imbalance found in the real building stock is exacerbated in the labelled data, increasing the overrepresentation of SFHs and THs, as opposed to ABs.

Like the size classes, the construction year classes also differ between census and TABULA. To address this mismatch, a uniform distribution of construction years within the respective year class was assumed and an exact year was assigned randomly to each building. This allows a flexible reaggregation of the required construction year classes.

The labelled dataset for construction year class classification includes approximately 6.9 million samples, encompassing 42 features (see Table 2). However, the distribution of the dataset is considerably skewed, as the 1949–1978 class makes up more than half (59.82 %) of the data (see Fig. 3). Furthermore, the earliest three classes together contain more than 80 % of the building samples, indicating a significant imbalance. As all labelled buildings are included in the model training process without additional sampling steps, this is the combined effect of the underlying distribution of the ground-truth data and the data labelling process described above. The imbalance in the training data can also be seen in the official statistics [47]. However, an over-representation of the majority class in the training dataset compared to the real building stock characteristics is apparent.

Class imbalance may lead to overfitting of the trained model on the dominating class, resulting in a loss of generalizability during inference, and so sample weights will be used during training as further outlined in Section 3.2.2.

The census data was utilized not only for extracting target labels for the size class and construction year classification but also to generate features. The number of inhabitants per grid cell was extracted and used as a feature, as well as the average number of occupants per household per grid cell, which was calculated by dividing the inhabitants per grid cell by the number of households per grid cell as reported by the census. This data was equally filtered for non-anonymized cells to exclude the possibility of randomly generated data being introduced into the dataset

**Table 1**
Mapping between census and TABULA size classes [45].

| TABULA size class | Census size classes |
| --- | --- |
| SFH | Detached single-family house |
| | Detached two-family house |
| | Single-family semi-detached house |
| | Two-family semi-detached house |
| TH | Single-family terraced house |
| | Two-family terraced house |
| MFH | Multi-family house with 3–12 apartments |
| AB | Multi-family house with 13 or more apartments |

**Table 2**
Number of occurrences in the labelled dataset for each of the two prediction targets.

| Prediction target | Number of samples | Number of features |
| --- | --- | --- |
| Size class | 10,813,067 | 41 |
| Construction year | 6861,243 | 42 |

and thus ensure a high degree of data quality.

The socio-economic dataset comprises information from eleven distinct sources. After cleaning the data, a bottom-up data integration strategy was employed, wherein data at identical regional levels were merged. This resulted in the formation of two datasets: one specific to the Local Administrative Unit (LAU) level and another corresponding to the Nomenclature of Territorial Units for Statistics (NUTS) level 3, due to the two resolution levels of the input dataset. These two datasets were merged by assigning the values from the NUTS-3 level to all LAU regions located within it. This process culminated in a consolidated dataset on the LAU level. This data was then added to the buildings as additional features by a spatial join.

Building and neighborhood morphological characteristics were extracted from the individual building data presented in Section 3.1.3 and used as features. Table 3 provides an overview of all features included in the three groups "building morphology," "neighborhood morphology," and "socio-economic characteristics."

### 3.2.2. Model training

In this study, the XGBoost machine learning algorithm was chosen on the basis of its numerous benefits over alternative decision tree classifiers such as Random Forest. A primary advantage of XGBoost is its capacity to process categorical variables, thereby avoiding the necessity for one-hot encoding, a method for transforming categorical variables into numerical data by expressing the categories in the form of binary vectors. This characteristic substantially reduces computational time. XGBoost processes categorical variables by considering each feature as a subgroup and executing a partition-based split on the condition that the value is an element of the categories [48]. This describes the fact that the condition for splitting is based on the membership of a value in a set of categories, where categories is a subset of all possible categories. Moreover, the datasets in this study contained numerous missing values across all features. Consequently, techniques like MICE imputation, a method that, when applied to a dataset with multivariate missing data, iteratively predicts the missing values for all features in the dataset through a combination of statistical assumptions and machine learning, were contemplated. However, upon testing, it was discerned that XGBoost's "Sparsity-aware Split Finding" technique for the efficient handling of missing values performed comparably, if not superiorly, and significantly reduced the computational load. This technique amasses statistics of non-missing entries in order to determine the optimal split for a value. In cases of missing values, the split was performed in the direction that yields the largest decrease in error so far. Statistics regarding the number of missing features per building are included in Supplementary Material S4.

Additionally, all target labels in this investigation, regardless of the classification problem, were heavily imbalanced. Although an approach akin to Wurm et al. [17] and Bandam et al. [32] utilizing SMOTE [49] was initially contemplated, it was abandoned due to the substantial increase in training time required, favoring sample weights instead. These sample weights can be provided directly to XGBoost when fitting the model. The weights were calculated using the sample method "balanced" from the sklearn library, which creates weights that are inversely proportional to the number of occurrences per class [50], thus giving higher weight to underrepresented classes. During training, the gradient of each sample and the loss function were computed and updated after each iteration. The sample weights affect how the model's parameters are updated. If a sample has a higher weight, its gradient will
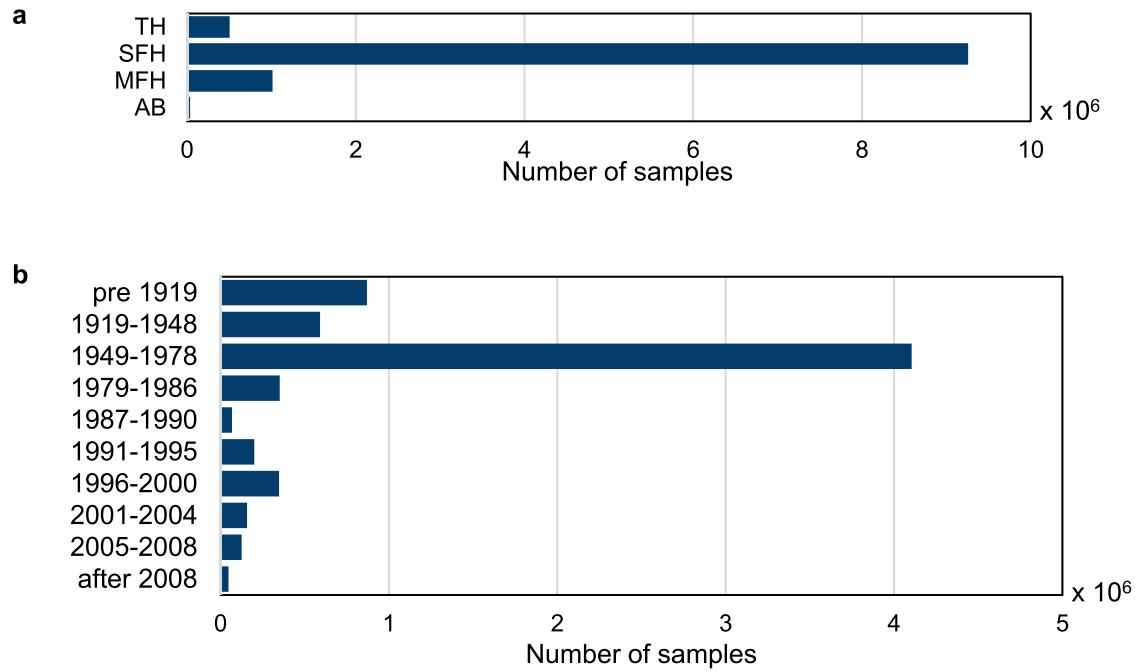
**Fig. 3.** Distribution of occurrences in the labelled datasets for: (a) size classes; and (b) construction years.

have a more substantial influence on the loss function and so a proportionately greater influence on the learning process. XGBoost ranked third out of 14 in a comparison study on multi-class imbalanced data classification methods and was the fastest multi-class boosting method in this study [51]. The first two, CatBoost and SMOTEBoost, could therefore be candidates for experiments in future studies. However, SMOTEBoost has training times of hours for datasets where CatBoost and XGBoost only take seconds [51], thus making it unsuitable for application on large datasets.

The size class and construction year models were both trained using a train-test-split of 80/20 alongside a five-fold cross-validation, employing stratified sampling. Furthermore, early stopping with four rounds was used to prevent overfitting. Both models were trained with a learning rate (*eta*) of 0.05, a maximum depth (*max_depth*) of 34, and a subsample ratio of 0.6 for constructing trees (*colsample_bytree*). Beyond these hyperparameters, default settings were applied. The hyperparameters were determined using a subset of the final dataset to diminish the training time and subsequently adjusted until performance progression plateaued. The hyperparameter progression used during testing is included in Supplementary Material S4.

### 3.2.3. Model evaluation

The models were evaluated using the F1-score. The F1-score (see Eq. (3)) combines the precision (see Eq. (1)) and recall (see Eq. (2)) metrics per class, and the weighted average F1-score is the average of the class-wise F1-scores weighted by the number of occurrences of each class. Precision measures the accuracy of the positive predictions made by the model and recall measures the ability of the model to identify all relevant instances within a dataset.

$$precision = \frac{true\ positives}{true\ positives + false\ positives} \tag{1}$$

$$recall = \frac{true\ positives}{true\ positives + false\ negatives} \tag{2}$$

$$F1 - score = \frac{2\ x\ precision\ x\ recall}{precision + recall} \tag{3}$$

In order to determine the performance of the model using the F1-score, scores are compared with model results from the literature (see Section 5) and against baseline models, which follow a simple assignment process. The baseline performances were calculated using the DummyClassifier from the scikit-learn library [52], using the strategies "most_frequent", which assigns the most frequently occurring category in the training data, and "uniform", which randomly assigns one of the occurring categories with equal probability, when predicting the output.

### 3.3. TABULA type assignment and validation

As described in Section 1, it is possible to determine the TABULA building code of each residential building based on the characteristics previously assigned (size class, construction year, refurbishment state). Within the TABULA typology, certain combinations of size classes and construction years are undefined, such as a TH built before 1960 (see Supplementary Material S5). However, the assignment of construction years and size classes is independent of these restrictions. To ensure internal consistency between those attributes within the final dataset, the assigning of inexistent TABULA types was not prevented. The extent of this inaccuracy is analyzed in the validation step. Furthermore, only non-regionally specified archetypes of the German TABULA typology were considered. Those that only are valid for the new federal states in Germany were disregarded.

For the validation, the construction year and size class distributions were compared to official statistics. This validation was carried out at the federal state level. Additionally, the accuracy of the assignment of TABULA types was assessed by comparing the total and relative deviation of TABULA type occurrences assigned by the methodology presented in this study with the building stock statistics provided in Loga et al. [45]. Due to the aggregation level of the statistics, this was only possible at the national level. It also reveals how many buildings were assigned an undefined TABULA type.

### 4. Results

The following sections present the performance of the trained machine learning models at the individual building level (see Section 4.1)

**Table 3**

Features for construction year and size class model. For some of the more complex building morphology features, the table also includes the calculation formula.

| Building morphology | |
| --- | --- |
| Perimeter | $P$ |
| Conditioned living area | $A\_C$ |
| Footprint area | $A$ |
| Roof type | |
| Average roof tilt | $t_{roof,\ avg}$ |
| Average roof height | $h_{roof.avg}$ |
| Roof area | $A_{roof}$ |
| Height | $h - h_{roof}$ |
| Average width of footprint polygon | $w_{avg} = (P/\pi * A)/(P^{(2/4)} * \pi)$ |
| Length | $l = A/w_{avg}$ |
| Wall area | $A_w = P * h$ |
| Surface area | $A_s = A + A_w + A_{roof}$ |
| Volume | $V$ |
| Sphericity | $S = \left(\pi^{1/3} * (6V)^{2/3}\right)/A_s$ |
| Normalized perimeter index | $(2 * \sqrt{\pi * A})/P$ |
| Diameter | $d$ |
| Proximity index: average distance between a footprint's centroid and all vertices | |
| Height-length-ratio | $l/h$ |
| *Size class | |
| Neighborhood morphology | |
| Grid complexity | |
| Number of touching residential buildings | |
| Length of shared walls | |
| Area of shared walls | |
| Residential buildings in [30,50, 100, 500] m radius | |
| Socio-economic characteristics | |
| Median income | |
| GDP per capita | |
| GDP per employee | |
| Unemployment rate | |
| Long-term unemployment rate | |
| Broadband expansion | |
| Population potential | |
| Employees in knowledge-intensive industries | |
| Academic employment rate | |
| Population with [high, middle, low] level of education | |
| Inhabitants per grid cell | |
| Average number of occupants per household | |
| RegioStaR17 region | |

*Only included in the construction year training dataset.

and the characteristics and plausibility checks of the resulting dataset aggregated to the federal state level (see Section 4.2).

### 4.1. Model performance

The model performances of the size class and construction year model were assessed both with regard to the overall performance and the results in the individual classes. Furthermore, feature importances per model were analyzed. The suitability of the selected features for training the machine learning models and the feasibility of the chosen model for predicting size class, and construction year of individual buildings were evaluated.
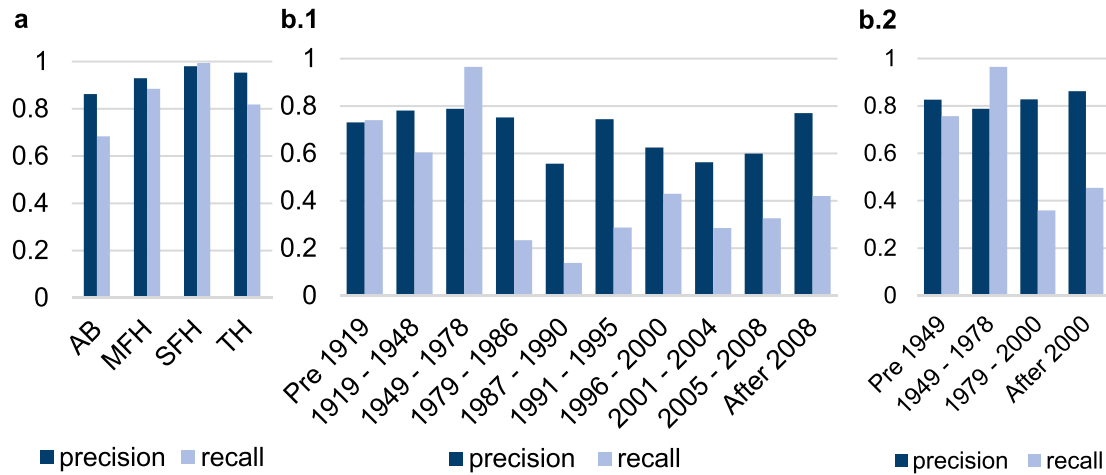
#### 4.1.1. Size class

The model designed for predicting the building size class achieved a weighted average F1-score of 97.4 % across five-folds for the test dataset. Machine learning models tend to perform well in the majority class but poorly in the minority ones, leading to biased predictions. Despite the data imbalance depicted in Fig. 3, precision (97.41 %) and recall (97.46 %) for the test dataset were almost equal, proving the effectiveness of the approach taken in applying sample weights. The model exceeded the best baseline model, which uses the "most-frequent" strategy, by 18.4 percentage points (see Supplementary Material S4).

Examining precision and recall for each class (see Fig. 4a), it can be observed that they surpass 90 % and 80 %, respectively, for all classes apart from ABs. Intriguingly, the precision score for SFHs falls below its recall, suggesting that while almost all SFHs in the dataset were correctly identified, other building classes were incorrectly classified as SFHs. Considering the scores for SFHs and THs together, the reduced recall score for THs could be attributed to them being incorrectly classified as SFHs. This in turn contributes, in conjunction with MFHs misclassified as SFHs, to the lower precision score for SFHs (see Supplementary Material S4). The low recall for ABs combined with a comparatively high precision indicates that while most buildings classified as ABs are in fact ABs, a relatively large share of ABs is not identified as such but could be misclassified as MFHs. One reason for this might be that the difference between ABs and MFHs lies purely in the number of housing units, as shown in Table 1. MFHs comprise buildings with 3–12 housing units, with ABs starting from 13 housing units upwards. The large range of housing units of MFHs likely results in a high diversity within this group. Combined with the seemingly arbitrary threshold of 13 housing units for ABs, this probably makes it difficult for the machine learning model to differentiate between the two size classes.
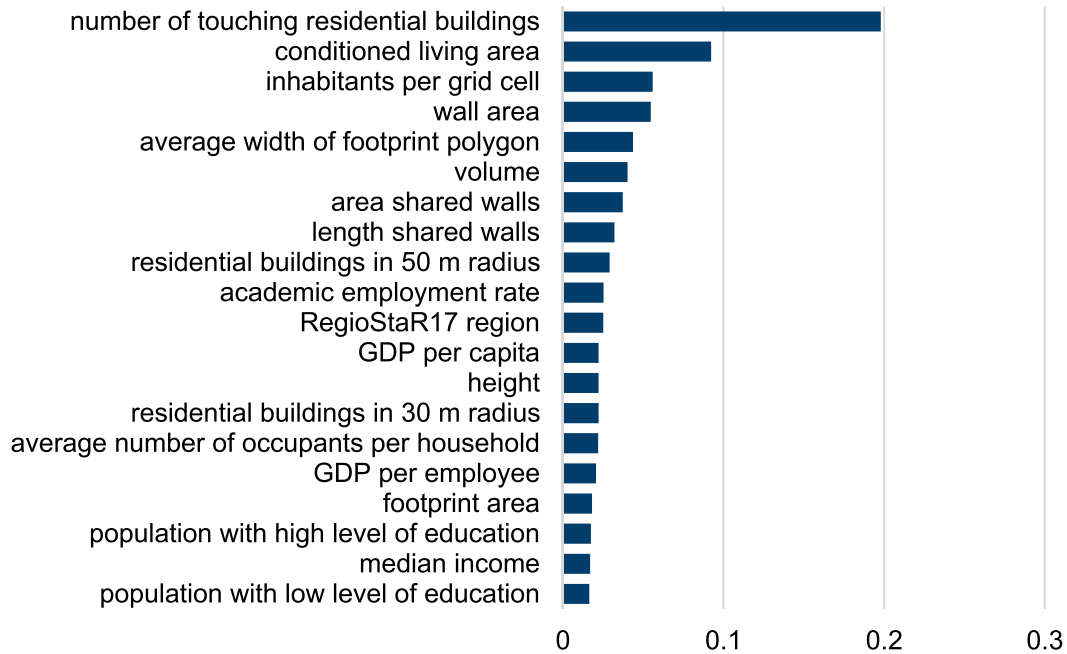
The ten most important features associated with building size class prediction (Fig. 5) provide information on the neighborhood morphology (number of touching buildings, the area of shared walls, the length of shared walls, and residential buildings in the 50 m radius), the individual building morphology (conditioned living area, wall area, average width, volume), and socio-economic characteristics (inhabitants per grid cell, academic employment rate). Further socio-economic features appear further down in the ranking. Some individual and neighborhood characteristics also appear to have a low degree of importance for size classification. This might be due in part to correlation between features, and features such as the footprint area and height being included in the conditioned living area, wall area, and volume. Whereas including correlated features is likely to influence feature importances, the model performance of decision tree algorithms, such as XGBoost, is expected to be relatively robust [53]. Therefore, these correlated features were all included in the model and not removed before training.

#### 4.1.2. Construction year

The model designed for predicting construction year classes achieved a weighted average F1-score of 73.93 % across five-folds for the test dataset. As previously mentioned, this metric combines precision and recall information per class and averages them across classes, taking the number of occurrences of the respective class into consideration. The model outperforms the baseline model using the "most-frequent" strategy by 6.33 percentage points (see Supplementary Material S4). The evaluation of the precision and recall for each class in the test dataset (see Fig. 4b.1) underscores the influence of the data imbalance on accuracy. Interestingly, the classes representing buildings constructed before 1979 exhibit the highest accuracy scores, which, if compared to the sample distribution in Fig. 3, implies a correlation between sample size and accuracy. Recall of the 1949–1978 class is very high, exceeding the precision, which indicates that while almost all buildings from this construction year period were identified as such, buildings from other construction year periods were mistakenly also assigned this label. The recall for underrepresented classes, notably the years between 1979 and 1995, is low, which suggests that the imbalance in the training dataset cannot be fully mitigated. However, the step width of the construction year classes is unbalanced, with a single class comprising between 3 and 30 years. Therefore, precision and recall were also calculated for regrouped classes (see Fig. 4b.2) with a more equal distribution of years. This shows a balanced precision between 78.8 % and 86.2 %. Recall on the other hand remains lower for the classes after 1979, which, in combination with the high recall of the 1949–1978 class, indicates that buildings from these classes were assigned earlier construction years

**Fig. 4.** Class-wise precision and recall on the test dataset for: (a) size class; (b.1) construction year period classification; (b.2) aggregated construction year period classification.



**Fig. 5.** 20 most important features for the building size class model.

than they actually had.

Reviewing feature importance, it can be identified that almost half of the features demonstrate an importance exceeding 2 % (see Fig. 6). Amongst the top ten features, socio-economic, individual building morphology, and neighborhood morphology features are relevant. Interestingly, building morphology features, while also positively contributing to the model's results and roof shape being among the top five features, are less important than for the size class model. This could indicate that the construction year is a characteristic shared by buildings in the same area and that the construction year period of a district manifests itself in the neighborhood morphology more than at the individual building level and that socio-economic features are more strongly correlated with the construction year period than the size class. Adding additional individual building features that are likely closely linked to the construction year period, such as façade configuration, are

expected to lead to a further improvement of the model.

### 4.2. Aggregated dataset characteristics and plausibility check

The final generated dataset including size class, construction year, refurbishment state, and TABULA type was analyzed regarding the distribution of classes. Furthermore, it was compared against official statistics for size class and construction year and against the TABULA report for TABULA type.

#### 4.2.1. Size class
Fig. 7 depicts the share of size classes of residential buildings for each of the 16 German federal states (NUTS-1 level) in the generated and training dataset and according to official statistics [54]. Overall, SFHs are the dominant size class. In most states, their share lies well above 50
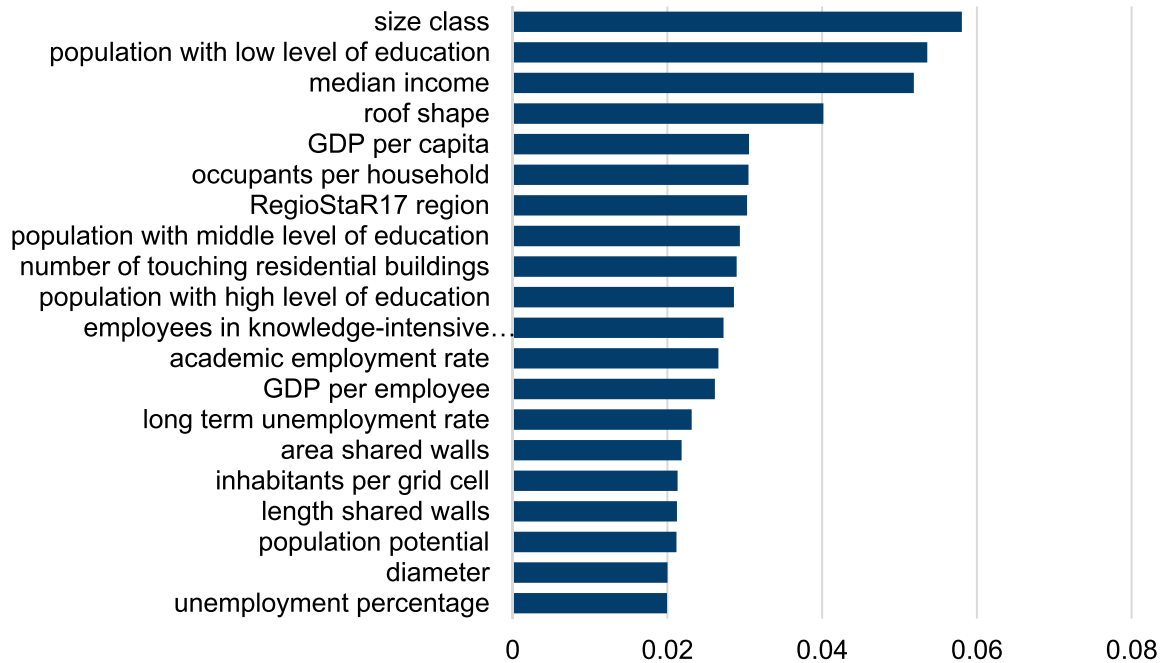
**Fig. 6.** 20 most important features for the construction year class model.

%, reaching almost 90 % in Brandenburg. The lowest shares of SFHs can be observed in the three city states of Bremen, Hamburg, and Berlin. ABs only constitute a significant share of buildings in Berlin, exceeding 10 %. A visual comparison shows that the distribution of size classes in the resulting dataset is similar to official statistics. There is a tendency that the share of SFHs and THs is overestimated in the generated dataset, whereas the share of ABs is underestimated. In all states, the share of SFH+TH is closer to the official dataset in the result than in the training data. On the other hand, apart from Berlin, where the share of ABs is highest according to the official data, the share of ABs in the resulting dataset is even lower than in the training dataset, mirroring the comparatively low recall for this class as reported in Section 4.1.1.

*4.2.2. Construction year*

Fig. 8 shows the average construction year of residential buildings at the NUTS-0 and NUTS-3 levels in the generated dataset. A tendency for older buildings to be located in the northeastern German states can be observed. Buildings with more recent construction years are found predominantly in the west and southeast. This pattern closely follows the former division of Germany into West and East Germany and is in line with the official statistics shown in Fig. 7, which highlight a large share of buildings in the "pre-1919″ and "1919–1949″ periods for Thuringia, Saxony, and Saxony–Anhalt, followed by Mecklenburg–Western Pomerania, Berlin, and Brandenburg. Although individual buildings can be of a much more recent construction year, the average for any NUTS-3 region is not beyond 1970. The NUTS-3 region with the youngest building stock is Vechta in Lower Saxony, with an average building construction year of 1969.

It can be observed that there are more buildings from the 1949–1979 period in the generated dataset than in the official statistics (see Fig. 8 and Fig. 7). These buildings are therefore missing from the other construction periods, leading to an underrepresentation of more recent construction years. In any case, the presented methodology appears to favor those construction year classes that are dominant. This could be due to the model being influenced by the class imbalance, despite the application of weights.

Fig. 7 depicts the distribution of construction year classes in the result and training dataset and compares it to the official statistics for each of the German federal states [47]. As in Fig. 3, the dominance of the 1950–1979 period in the training dataset can be observed, as well as a high level of occurrence of the pre-1919 and 1919–1949 periods. Overall, the distribution pattern of the result mirrors that of the official federal state statistics. The tendency to overrepresent majority classes, which was already apparent for the training data, is also visible in the result dataset, however. This is especially pronounced for, e.g., Saarland and Rhineland–Palatinate, whereas the distribution is more similar in Saxony, Saxony–Anhalt, and Thuringia.
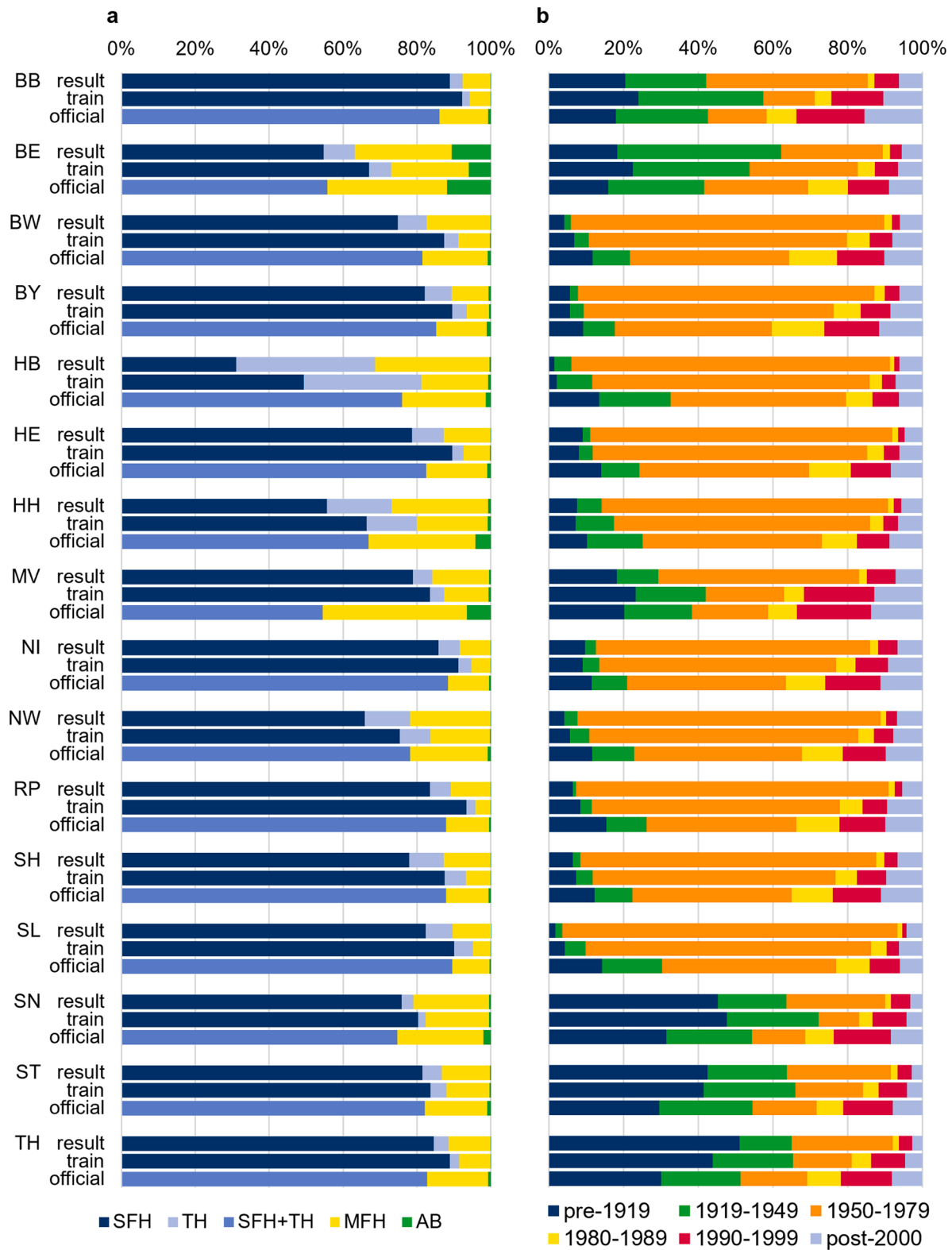
*4.2.3. Refurbishment state*

Fig. 9 shows the distribution of refurbishment states in the German federal states. This distribution corresponds to the data provided in Metzger et al. [28], which the randomized assignment was based upon. Most buildings have undergone some refurbishment, while buildings that have been fully refurbished and are in an advanced state of refurbishment constitute the smallest group. The highest share of at least partially refurbished buildings can be found in Saxony-Anhalt, Saxony, Thuringia, Brandenburg, and Mecklenburg-Western Pomerania, with 70–73 %, the lowest share in Hamburg with 50 %.

*4.2.4. TABULA type*

The total number of buildings per TABULA type in the generated dataset is shown in Fig. 10. As can also be seen in Fig. 7, SFHs are the predominant size class. Amongst the SFHs, there is a similar number from the second and fourth to the sixth generations, i.e., 1860–1918 and 1949–1978. The gap in the third generation, i.e., 1919–1948, aligns with the Second World War period (1939–1945) in Germany, which saw lower construction activity in the building sector. This drop, though not as pronounced as for SFHs, can be seen across all size classes. A sharp drop in the number of buildings constructed after 1978 can also be observed.
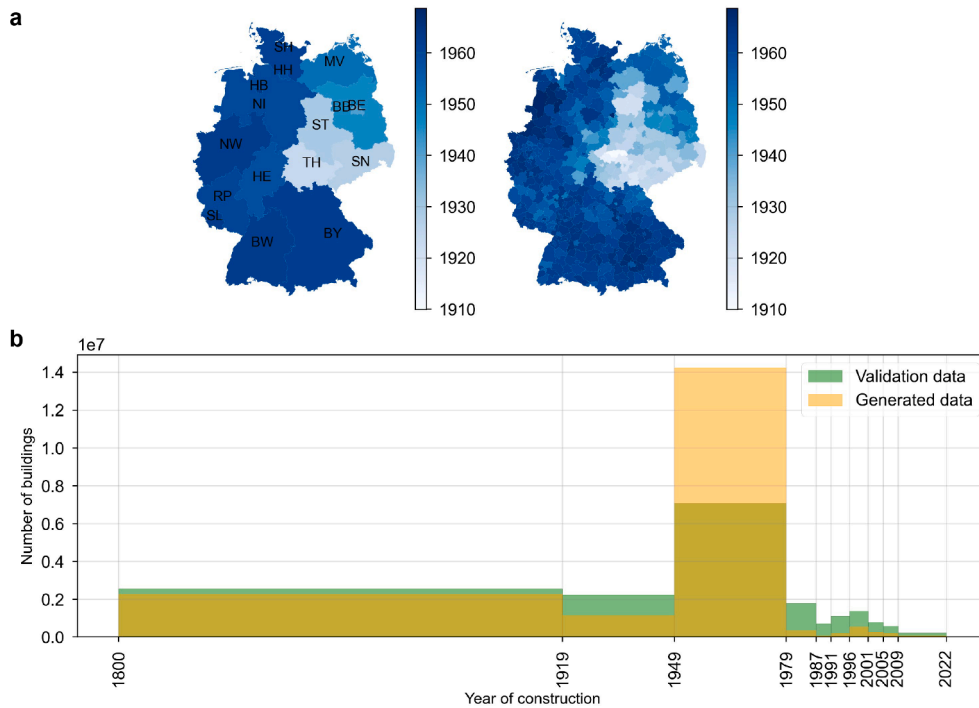
The influence of the overestimation of buildings in the construction year classes mentioned above can also be observed in Fig. 11, which displays the difference between the areas of TABULA types (neglecting the refurbishment state) in the generated dataset and those according to the TABULA report [45]. Values are combined into 5 clusters based on a KMeans cluster analysis of the respective TABULA type's specific heat
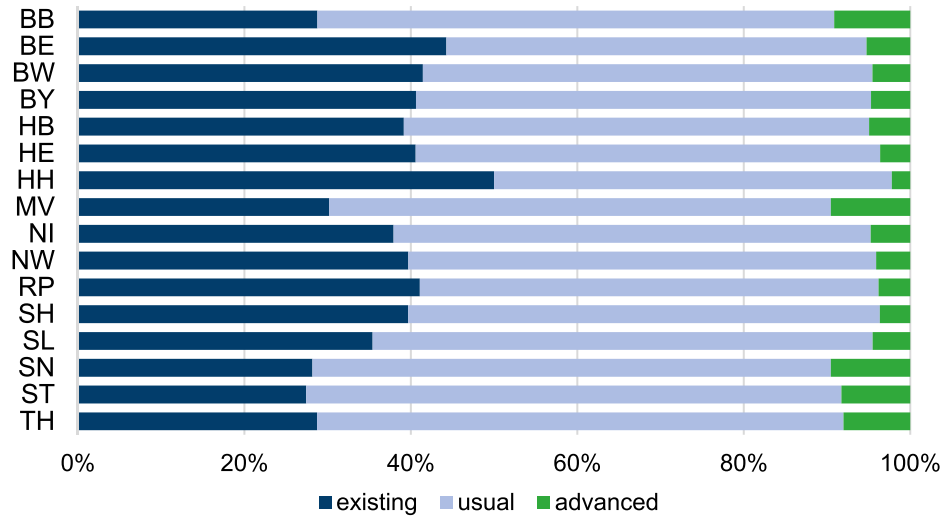
**Fig. 7.** Share of (a) size classes and (b) construction years in the German federal states in training and the resulting dataset and official statistics [54]. In the official statistics, the most recent category for construction years is 2000–2011.

demand. Cluster 1 comprises the 1950–1979 construction year class for SFHs and partly that for THs. As discussed above, the construction year class is dominant and a tendency to ascribe this class to buildings has been observed. Combined with the previously discussed over-representation of SFHs, this leads to the high deviation of the reference

building area between the TABULA report and the generated dataset. Cluster 1 and 3 contain the more recent construction year classes, which were shown to be more difficult to detect for the model, resulting in lower building areas in the generated dataset than in the TABULA report. This illustrates the complexity of the assignment problem at

**Fig. 8.** Distribution of construction years in the generated dataset. (a) Average construction year of residential buildings at the NUTS-0 and NUTS-3 levels; (b) comparison with official statistics for Germany [55].



**Fig. 9.** Distribution of refurbishment states in the German federal states, assigned based on Metzger et al. [28].

hand, as it relies on multiple previous processing steps.

In order to address the aforementioned issue and provide a more meaningful validation that abstracts from the concrete building types and focuses more on the applicability of the results dataset, the heat demand calculated based on the TABULA types was considered (see Table 4). Table 4 displays the heat demand calculated by multiplying: a) the areas of TABULA types in the generated dataset considering only size class and construction year class; and b) the areas according to the TABULA report, with the specific heat demand of the respective building type according to TABULA for all three refurbishment states. The heat demand based on the generated dataset is higher than that based on the TABULA report. When also considering the refurbishment state of individual buildings in the generated dataset, this results in a heat demand

of 466 TWh, which is 2.4 % higher than the value of 455 TWh for space heating reported by the German Environment Agency for 2022 [56] and 7.4 % lower than the value for 2020 reported by the German Federal Office for Statistics [57]. Even when assuming that all buildings are in the lowest refurbishment state, the heat demand calculated based on the TABULA report alone remains 2.9 % below the value reported by the German Environment Agency [56]. This can in part be attributed to the fact that the TABULA report dates back to 2012 and is therefore outdated with regard to building numbers. This indicates that although the numbers of buildings by TABULA type in the generated dataset deviate from the TABULA report, the generated dataset is nevertheless useful for estimating the current heat demand.
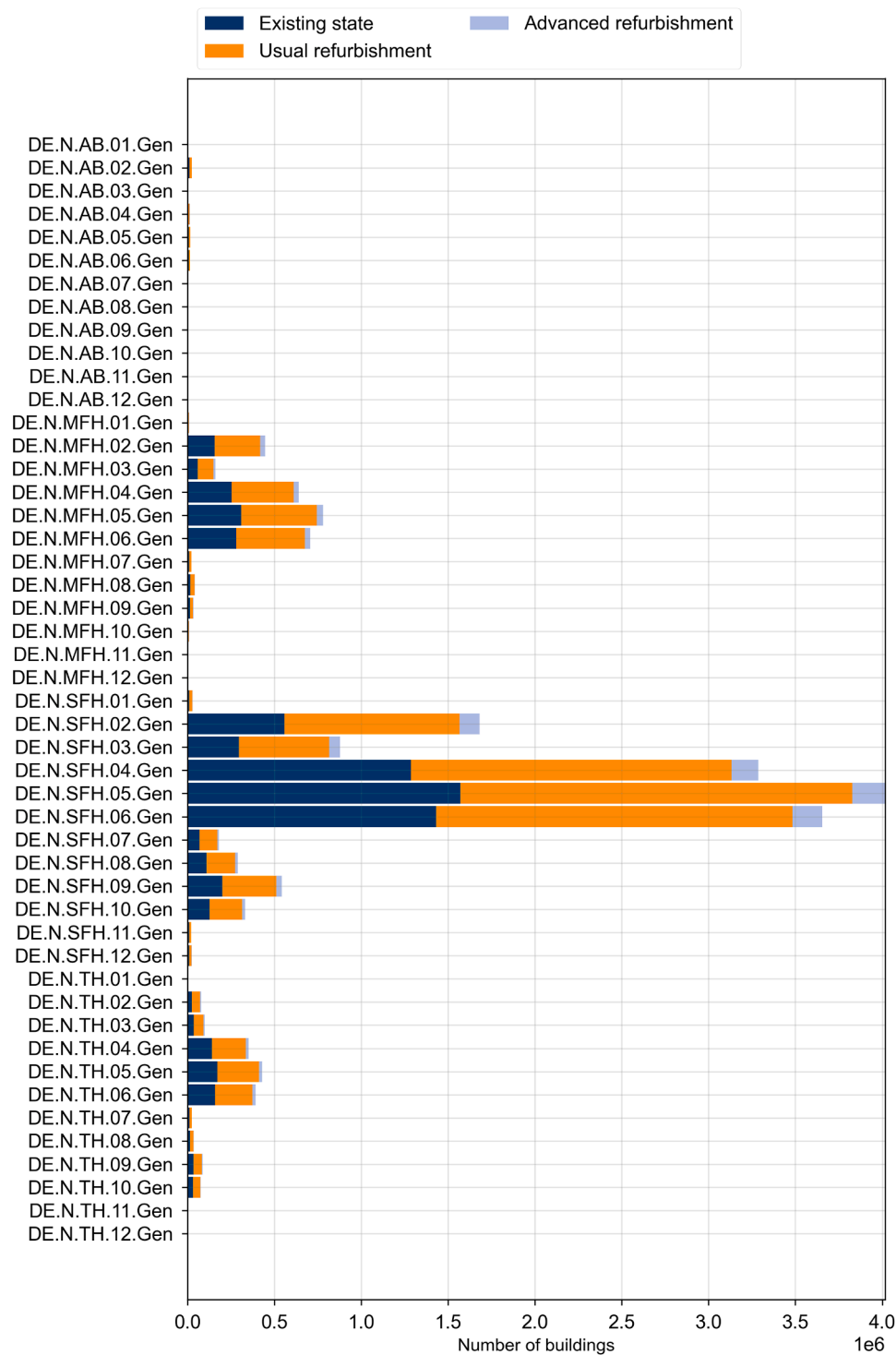
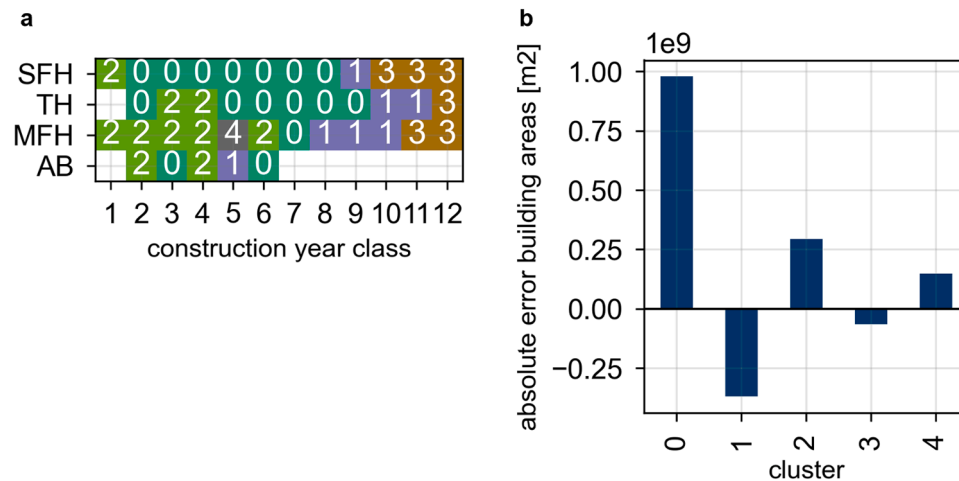**Fig. 10.** Number of TABULA types in the generated dataset.

## 5. Discussion

This study shows that assigning TABULA types to all residential buildings using a machine learning model trained with data from the German census is feasible. The performance of the machine learning model showed that it is possible to predict size classes with a high degree of accuracy. The construction year prediction is more challenging, which is at least partly due to the higher number of target classes. Both models clearly outperform the baseline model. Both machine learning tasks had unbalanced training data. The results indicate that using

sample weights was useful in balancing performance across classes. However, there were still some imbalances in both the size class and construction year prediction performance.

With a weighted average F1-score of 97.4 %, the size class model outperformed those by Wurm et al. [17] and Droin et al. [18] who follow a similar approach and reached scores of 96 % and 95 %, respectively. The construction year class model performance of 73.93 % falls below what was reported by Rosser et al. [20] and Garbasevschi et al. [19], who report scores of 77 % and 80 %, respectively. However, the former only consider five construction year periods, half as many as in this

**a**



**b**



**Fig. 11.** (a) TABULA types clustered by specific heat demand and (b) absolute deviation between building area per TABULA type cluster in the generated dataset and the report of the TABULA project.

**Table 4**

Comparison between heat demand for space heating reported by the German Environment Agency [56], calculated based on the generated dataset including the assigned refurbishment states, and extrapolated from the TABULA report and the generated dataset assuming that all buildings are in only one of the three refurbishment states: 'existing', 'usual', or 'advanced'.

| Source | Total heat demand [TWh] |
| --- | --- |
| ETHOS.BUILDA | 466 |
| Environment Agency UBA [56] | 455 |
| German Federal Office for Statistics [57] | 503 |
| TABULA report (all 'existing' state) | 442 |
| ETHOS.BUILDA (all 'existing' state) | 614 |
| TABULA report (all 'usual' state) | 286 |
| ETHOS.BUILDA (all 'usual' state) | 383 |
| TABULA report (all 'advanced' state) | 200 |
| ETHOS.BUILDA (all 'advanced' state) | 275 |

study, and only train their model for one city. The latter predict the same number of construction year periods as in this study but also limit their study to a few cities, all located in one federal state and thus in relatively close spatial proximity. Therefore, both probably deal with a more homogeneous building stock, which makes it easier for the model to learn the connection between characteristics and construction year periods. However, this could significantly reduce the amount of available training data per model, which in turn reduces its ability to learn and could therefore also lead to reduced accuracies. Considering the heterogeneity of all German building stock, the performance of the construction year model in this study is therefore competitive. In the future, adding more features, such as façade structures and material, could be tested to improve the model performance.

Despite the high performance of the models, some differences between the generated dataset and official statistics of German building stock at the federal state level can be observed. With respect to the construction year, the share of buildings assigned to the 1949–1979 period is too high in the generated data compared to official statistics. This issue is caused by a combination of effects that propagate and magnify through the processing workflow. First, the aforementioned construction year period is in fact the dominant period for buildings in Germany. Second, during the selection of census grid cells, only those cells that include one manifestation of the construction year period for all buildings within a cell are selected for the training data generation. This favors the majority class, as it is less likely that homogenous grid cells with only one of the rarer construction year periods exist. This amplifies the imbalance present in the training dataset. Third, when training a machine learning model on this training data, despite using

weights, the model tends to favor the majority class. Using this model then exacerbates the imbalance in the final dataset. A similar explanation is also applicable to size classes and explains the overrepresentation of SFHs in the generated dataset. In sum, this leads to an overly large proportion of SFHs in the 1949–1979 period. In order to mitigate this issue, one could test a less restrictive creation procedure for the training dataset from the census data with the goal of including more of the minority classes by not only considering homogeneous grid cells. This could, for example, mean that a grid cell with mainly ABs will also be considered and the 'AB' size class would be assigned to the largest buildings within the grid cell. Assessing the availability and potentially including additional datasets for adding new features or labelling more buildings for training, such as real estate data or local construction year datasets, might also be options for further studies. Furthermore, it should be kept in mind that discrepancies, for example in building numbers, can also be caused by characteristics of the underlying base data, especially OpenStreetMap data, as discussed in Dabrock et al. [6] and Bandam et al. [32]. Finally, other approaches to handling data imbalances should be tested, such as random under-sampling.

Assigning the TABULA type to buildings based on construction year, size class and refurbishment state is straightforward once these three characteristics are successfully assigned. After a TABULA type is assigned, this gives access to a range of typical properties of the respective building archetype defined by the TABULA project, such as U-values of building components or typical specific heat demands. This provides a valuable basis for further analyses of the building sector, particularly in the context of energy analysis. However, as the TABULA type is directly derived from the three aforementioned attributes, the issues explained above that might lead to a misassignment of one of those also leads to the assignment of a wrong TABULA type. This is especially evident for the fourth to sixth generation of SFHs, as this is where the overrepresentation of the majority class in both construction year and size class coincides. Most likely, buildings are wrongly classified as a similar type, e.g., an SFH as a TH of the same generation. Unfortunately, there is no individual building data or even statistical data on the distribution of TABULA types below the federal state level available that could be used for validation and allow analyses regarding the misclassification of types and its impact on further analyses. However, the heat demand calculated based on the assigned TABULA types in the generated dataset has been shown to be close to the value reported by the Environment Agency. This proves the plausibility and applicability of the generated dataset.

A challenge in the decoupling of the TABULA type assignment from the construction year, size class, and refurbishment assignment is that it potentially results in the assignment of types that are not defined within

the TABULA archetype framework. The validation shows that this is the case for less than 10,000 buildings, i.e., only about 0.05 % of all residential buildings in Germany. The impact of handling those buildings and, for example, mapping them to a TABULA type that is assumed to be most similar in its characteristics, should be investigated in subsequent analyses.

Another limitation of the presented approach is the lack of data on the refurbishment state of buildings, especially on an individual building level. Therefore, the assignment had to rely on statistics at the level of federal states. In reality, there are probably smaller scale geographical variances in the share of refurbished buildings. For instance, differences between urban and rural areas or more affluent and poorer city districts may exist. This should be the object of further research. Furthermore, it should be noted that the current results are based on data from 2011 for the census data [23] and 2019 for the refurbishment data [28]. While it can be argued that construction years and size classes of existing buildings are constant and, as mentioned in the introduction, the new construction rate is low, it is likely that the share of refurbished buildings shows a relevant increase in the nearer future, thus requiring an update of the result data as soon as new refurbishment data becomes available.

The TABULA project provides building typologies for most of the EU countries. For these countries, the overall workflow presented in this study might be applied. However, the transferability to other regions depends strongly on data availability and quality. For example, for the Netherlands an open nationwide 2-D and 3-D building dataset including detailed information such as the construction year [58] and open census data at the neighborhood level including the share of single- and multi-family houses [59] are available, providing a good basis for the TABULA assignment. For Austria, on the other hand, 3-D building data is only available for some regions, such as Styria [60] and Tyrol [61], and the census data with a spatial resolution of 250 m is expensive and may only be used internally [62], which makes the adaptation of the workflow more challenging. Furthermore, there is a great variability of OpenStreetMap data quality across regions [63]. In regions where building level data on construction year, size class, or refurbishment state exists, parts of the workflow can be simplified. Instead of training a machine learning model, simple direct assignment or mapping procedures might be sufficient. On the other hand, regions with even fewer datasets available than in Germany could pose a problem in terms of training data generation. If no nationwide datasets similar to the census are available, alternative approaches should be considered. One approach could be gathering sample data that enables the deduction of target labels for training a machine learning model. Alternatively, models trained in one geographic region could potentially also be applied to other regions, given a similar morphological structure of the building stock. It would, for example, be interesting to test the model in other countries in Central Europe. This would still require regional data, at least for validation, however. Additionally, the training of the model for Germany relies on socio-economic datasets that are likely not available in the same format in other countries. Future studies could explore how the general approach, which is transferable to other regions and has been proven successful for Germany, could be adapted to local data availability.

The resulting open access dataset contains not only TABULA types for all residential buildings in Germany but also the constituting attributes, namely size class, construction year and refurbishment state, as well as some of the basic building morphological features, such as footprints and heights. The dataset is valuable for further research in the field of building energy demand analysis, for developing targeted decarbonization pathways, or for any kind of building-related analysis requiring spatial data.

When developing decarbonization pathways for the German building stock, analyses are often top-down and based on statistical data, e.g., Thomas et al. [64], Prognos et al. [65] and Fraunhofer IWES/IBP [66]. However, spatially resolved building level data offers many advantages

and is the basis for bottom-up studies. For example, the data can be used to configure models for individual buildings, such as ETHOS.HiSim [14], as shown in Rieck et al. [15]. At the same time, the TABULA types combined with additional building level data provide an easy way to cluster the individual buildings, following a combination of the building-by-building approach described by Mastrucci et al. [67] and the archetype approach described by Swan et al. [68], in order to draw conclusions for the entire building stock while minimizing the computational resources required for the simulation. In contrast to the aforementioned top-down studies, the dataset created in this study, in conjunction with additional geospatial data, also allows the inclusion of spatial constraints and potentials relevant to decarbonization options, such as the availability of space for heat pumps [69,70], the proximity to clean heat sources for district heating [71], or small-scale regional geothermal potential [72]. Thus, the dataset allows to define feasible decarbonization options for individual buildings and can be the basis for detailed decarbonization pathways beyond national statistical analyses.

However, it should be noted that while the dataset contains attributes for individual buildings, only the size class and construction year are validated against a test set at the individual building level. Due to the lack of refurbishment data at the individual building level, no statements can be made at this time regarding the accuracy of refurbishment states and TABULA types at an aggregation level below the federal state. Therefore, simulations based on this data should not be considered accurate for each individual building. Instead, results should be aggregated in order to draw reliable conclusions. The data is available through an API[1] and can be easily accessed via a Python client.[2] This mode of access combined with filtering options facilitates the integration of building data into existing Python-based workflows and removes the necessity to download and handle unnecessarily large and complex datasets. Therefore, this dataset has the potential to significantly simplify and speed up data-intensive research tasks and thereby contributes to efforts in the research community to advance the energy transition.

## 6. Conclusion and outlook

This study demonstrated the feasibility of a methodology for assigning TABULA types to all residential buildings in Germany by leveraging machine learning. The census data, in combination with OpenStreetMap and additional datasets, are suitable for deriving training data for the machine learning task. The extensive feature set, which includes building and neighborhood morphological features as well as socio-economic ones, enables predictions of size and construction year classes with a high level of performance, with F1-scores of 97.4 % and 73.93 %, respectively. This shows that training and applying a nationwide machine learning model is possible, thereby going beyond the small spatial scope previously seen in the literature. The lack of proper data on refurbishment states constitutes a limitation, and therefore the acquisition of building level refurbishment data could significantly enhance the reliability of the results. Additionally, a tendency to overestimate the occurrence of majority classes, namely SFHs in the 1949–1979 construction year period, which also has an effect on the TABULA type distribution in the resulting dataset, has been observed and discussed. However, the methodology is well able to capture the shares and spatial distribution of attributes within the German building stock.

The published open data, including TABULA types, size class, construction year, refurbishment state at the individual building level is a valuable asset for the research community. It helps to fill the gap in building-related data.

In subsequent studies, analyzing the transferability of the workflow

---

to regions outside of Germany would be interesting. The effort to adapt the workflow depends on the data availability. Due to the good data availability and the close spatial proximity suggesting a similarity in the building stock, applying the workflow to regions in the Netherlands would be a reasonable and interesting next step. Finally, this study focuses exclusively on residential buildings. It would be interesting to develop a similar methodology for non-residential buildings, for example using the archetypes developed by Hörner and Bischof [73].

### Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work, the authors used ChatGPT for translations and linguistic revisions in Sections 3.1.1, 3.1.4, 3.2.1, 3.2.2, and 4.1. Furthermore, DeepL was used for linguistic revisions in Sections 5 and 6. After using these tools, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

### CRediT authorship contribution statement

**Kristina Dabrock:** Writing – review & editing, Writing – original draft, Visualization, Validation, Supervision, Software, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Jens Ulken:** Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation. **Noah Pflugradt:** Writing – review & editing, Supervision, Methodology, Conceptualization. **Jann Michael Weinand:** Writing – review & editing, Supervision. **Detlef Stolten:** Supervision, Resources, Project administration, Funding acquisition.

### Declaration of competing interest

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

### Acknowledgements

### Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.buildenv.2025.112782.

### Data availability

https://doi.org/10.5281/zenodo.13771740 (The data is available at:)

### References

[1] European Commission, focus: Energy efficiency in buildings," European Commission - European Commission, 2022. Accessed: Apr. 04[Online]. Available: https://ec.europa.eu/info/news/focus-energy-efficiency-buildings-2020-lut-17_en.

[2] Bundesministerium für Wirtschaft und Klimaschutz (BMWi), "Energieeffizienz in Zahlen 2021 [Energy efficiency in figures 2021]," 2021.

[3] European Commission, "EU building stock observatory." Accessed: Mar. 22, 2022. [Online]. Available: https://energy.ec.europa.eu/topics/energy-efficiency/energy-efficient-buildings/eu-building-stock-observatory_en.

[4] Joint Research Centre, "Jrc-Idees." Accessed: May 23, 2023. [Online]. Available: https://joint-research-centre.ec.europa.eu/potencia/jrc-idees_en.

[5] co2online, "Gebäudedaten [building data]." Accessed: May 23, 2023. [Online]. Available: https://www.wohngebaeude.info/.

[6] K. Dabrock, N. Pflugradt, J.M. Weinand, D. Stolten, Leveraging machine learning to generate a unified and complete building height dataset for Germany, Energy AI 17 (2024) 100408, https://doi.org/10.1016/j.egyai.2024.100408. Sep.

[7] M. Schwanebeck, M. Krueger, R. Duttmann, Improving GIS-based heat demand modelling and mapping for residential buildings with census data sets at regional and sub-regional scales, Energies 14 (4) (2021) 1029, https://doi.org/10.3390/en14041029. Feb.

[8] X. Yang, et al., A combined GIS-archetype approach to model residential space heating energy: a case study for The Netherlands including validation, Appl. Energy (2020), https://doi.org/10.1016/j.apenergy.2020.115953.

[9] Tabula Project Team, Typology Approach for Building Stock Energy Assessment. Main Results of the TABULA Project, IWU Institut Wohnen und Umwelt, Darmstadt, 2012. Accessed: Jan. 03, 2023. [Online]. Available: https://episcope.eu/fileadmin/tabula/public/docs/report/TABULA_FinalReport.pdf.

[10] ArGe Arbeitsgemeinschaft für zeitgemäßes Bauen eV, Gebäudetypologie Schleswig–Holstein. Leitfaden für wirtschaftliche und energieeffiziente Sanierungen verschiedener Baualtersklassen [Building typology for Schleswig–Holstein. Guidelines for cost-effective and energy-efficient renovation of buildings of different ages], Bauen in Schleswig–Holstein 47 (2012).

[11] "TABULA WebTool." Accessed: Apr. 08, 2024. [Online]. Available: https://webtool.building-typology.eu/#bm.

[12] "Wohngebäude in Deutschland nach Baujahr [Residential buildings in Germany by year of construction]," Statista. Accessed: Jan. 03, 2025. [Online]. Available: https://de.statista.com/statistik/daten/studie/1385022/umfrage/wohngebaeude-in-deutschland-nach-baujahr/.

[13] D. Heidenthaler, M. Leeb, P. Reindl, L. Kranzl, T. Bednar, M. Moltinger, Building stock characteristics of residential buildings in Salzburg, Austria based on a structured analysis of energy performance certificates, Energy Build. 273 (2022) 112401, https://doi.org/10.1016/j.enbuild.2022.112401. Oct.

[14] N. Pflugradt, ETHOS.HiSim - House infrastructure simulator. (Dec. 30, 2024). Python. Accessed: Jan. 03, 2025. [Online]. Available: https://github.com/FZJ-IEK3-VSA/HiSim.

[15] K. Rieck, K. Dabrock, N. Pflugradt, J.M. Weinand, D. Stolten, Large-Scale Quantification of the Future Self-Covered Heat Demand Using a Nationwide Residential Building Database, Social Science Research Network, Rochester, NY, 2024 4916684. Aug. 05Accessed: Jan. 03, 2025. [Online]. Available: https://papers.ssrn.com/abstract=4916684.

[16] L.A. Blanco Bohorquez, M. Aditya, B. Schiricke, B. Hoffschmidt, Classification of Building Properties from the German Census Data for Energy Analysis Purposes, presented at the Building Simulation 2023, Shanghai, China, 2023. JunAccessed: Nov. 29, 2023. [Online]. Available: https://elib.dlr.de/199041/.

[17] M. Wurm, A. Droin, T. Stark, C. Geiss, W. Sulzer, H. Taubenboeck, Deep learning-based generation of building stock data from remote sensing for urban heat demand modeling, ISPRS Int. J. Geo-Inf. 10 (1) (2021) 23, https://doi.org/10.3390/ijgi10010023. Jan.

[18] A. Droin, M. Wurm, W. Sulzer, Semantic Labelling of Building Types, A comparison of two approaches using Random Forest and Deep Learning, 2020.

[19] O.M. Garbasevschi, et al., Spatial factors influencing building age prediction and implications for urban residential energy modelling, Comput. Environ. Urban Syst. 88 (2021) 101637, https://doi.org/10.1016/j.compenvurbsys.2021.101637. Jul.

[20] J.F. Rosser, D.S. Boyd, G. Long, S. Zakhary, Y. Mao, D. Robinson, Predicting residential building age from map data, Comput. Environ. Urban Syst. 73 (2019) 56–67, https://doi.org/10.1016/j.compenvurbsys.2018.08.004.

[21] S. Becker et al., "Metastudie zur verbesserung der datengrundlage im Gebäudebereich - Leistung gemäß rahmenvertrag zur beratung der abteilung II des BMWK [Meta-study to improve the data basis in the building sector - service according to the framework contract for advising department II of the federal ministry of economy and climate protection (BMWK)]," Berlin, München, 2022. Accessed: Nov. 29, 2023. [Online]. Available: https://www.bmwk.de/Redaktion/DE/Publikationen/Energie/metastudie-verbesserung-datengrundlage-gebaeudebereich.pdf?_blob=publicationFile&v=1.

[22] D.K. Alexander, S. Lannon, and O. Linovski, "The identification and analysis of regional building stock characteristics using map based data," p. 8, 2009.

[23] Statistische Ämter des Bundes und der Länder, "Zensus 2011 - Gebäude- und wohnungsbestand in Deutschland - endgültige ergebnisse [Census 2011 – building and housing stock in Germany – final results]." 2015. Accessed: May 22, 2023. [Online]. Available: https://www.zensus2011.de/SharedDocs/Downloads/DE/Publikationen/Aufsaetze_Archiv/2015_12_NI_GWZ_endgueltig.pdf?_blob=publicationFile&v=2.

[24] M. Zirak, V. Weiler, M. Hein, U. Eicker, Urban models enrichment for energy applications: challenges in energy simulation using different data sources for building age information, Energy 190 (2020) 116292, https://doi.org/10.1016/j.energy.2019.116292.

[25] M. Zeppelzauer, M. Despotovic, M. Sakeena, D. Koch, M. Döller, Automatic prediction of building age from photographs, in: Proceedings of the 2018 ACM On International Conference On Multimedia Retrieval, in ICMR '18, New York, NY, USA, Association for Computing Machinery, Jun. 2018, pp. 126–134, https://doi.org/10.1145/3206025.3206060.

[26] Y. Li, Y. Chen, A. Rajabifard, K. Khoshelham, and M. Aleksandrov, "Estimating building age from google street view images using deep learning (Short paper)," p. 40:1–40:7, 2018, doi: 10.4230/LIPIcs.GISCIENCE.2018.40.

[27] F. Biljecki, M. Sindram, Estimating building age with 3d gis. ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences, Copernicus

GmbH, Oct. 2017, pp. 17–24, https://doi.org/10.5194/isprs-annals-IV-4-W5-17-2017.

[28] S. Metzger, K. Jahnke, N. Walikewitz, M. Otto, A. Grondey, S. Fritz, Hintergrundbericht: Wohnen und Sanieren [Background report: Housing and Renovation, Umweltbundesamt, 2019 Accessed: Jul. 24, 2024. [Online]. Available: https://www.umweltbundesamt.de/publikationen/hintergrundbericht-wohnen-sanieren.

[29] "Zensus 2011 - Methoden und verfahren [2011 census – methods and procedures]," 2015.

[30] Statistische Ämter des Bundes und der Länder, "Ergebnisse des zensus 2011 zum download - erweitert [results of the 2011 census available to download – expanded]." 2020. Accessed: Dec. 19, 2022. [Online]. Available: https://www.zensus2011.de/DE/Home/Aktuelles/DemografischeGrunddaten.html.

[31] Geofabrik GmbH, "Geofabrik download server." Accessed: May 24, 2023. [Online]. Available: https://download.geofabrik.de/europe/germany.html.

[32] A. Bandam, E. Busari, C. Syranidou, J. Linssen, D. Stolten, Classification of building types in Germany: a data-driven modeling approach, Data 7 (4) (2022) 45, https://doi.org/10.3390/data7040045. Apr.

[33] "building | Keys | Values | OpenStreetMap Taginfo Germany." Accessed: Jan. 06, 2025. [Online]. Available: https://taginfo.geofabrik.de/europe:germany/keys/building#values.

[34] Arbeitsgemeinschaft der Vermessungsverwaltungen der Länder der Bundesrepublik Deutschland (AdV), "Dokumentation zur modellierung der geoinformationen des amtlichen vermessungswesens (GeoInfoDok) - ALKIS-objektartenkatalog [Documentation on the modeling of geospatial information for official cadastral surveying (GeoInfoDok) - ALKIS object type catalog]." Apr. 11, 2008. Accessed: May 24, 2023. [Online]. Available: https://www.adv-online.de/GeoInfoDok/GeoInfoDok-6.0/binarywriterservlet?imgUid=8f830072-8de8-9221-d5ad-8f138a438ad1&uBasVariant=11111111-1111-1111-1111-111111111111.

[35] SIG 3D, "CityGML code list roofType." 2012. Accessed: May 23, 2023. [Online]. Available: http://www.sig3d.org/codelists/citygml/2.0/building/2.0/_AbstractBuilding_roofType.xml.

[36] "Key:roof:shape – openstreetmap wiki." Accessed: May 23, 2023. [Online]. Available: https://wiki.openstreetmap.org/wiki/Key:roof:shape.

[37] F. Biljecki and H. Ledoux, Solar3Dcity. (May 03, 2023). Python. 3D geoinformation research group at tu delft. Accessed: May 23, 2023. [Online]. Available: https://github.com/tudelft3d/Solar3Dcity/blob/8c755fda852f234343f3ff27e1061dccfff84a2c/polygon3dmodule.py.

[38] "building | Keys | Combinations | OpenStreetMap Taginfo Germany." Accessed: Jan. 06, 2025. [Online]. Available: https://taginfo.geofabrik.de/europe:germany/keys/building#combinations.

[39] S. Pekka, V. Michael, Starke mietpreissteigerungen und erste aufwärtstendenzen bei wohnungspreisen [Sharp rent increases and first signs of an upward trend in apartment prices], Sagner IW-Report (6) (2024). FebAccessed: Apr. 08, 2024. [Online]. Available: https://www.iwkoeln.de/studien/pekka-sagner-michael-voigtlaender-starke-mietpreissteigerungen-und-erste-aufwaertstendenzen-bei-wohnungspreisen.html.

[40] "Schlechte energieeffizienz drückt die preise [Poor energy efficiency depresses prices]," ImmobilienScout24. Accessed: Apr. 08, 2024. [Online]. Available: https://www.immobilienscout24.de/wissen/verkaufen/energieeffizienz-und-preise.html.

[41] L. Edlund, C. Machado, and M. Sviatchi, "Bright minds, big rent: gentrification and the rising returns to skill".

[42] S. Maretzke, J. Ragnitz, G. Untiedt, Betrachtung Und Analyse Von Regionalindikatoren zur Vorbereitung des GRW-Fördergebietes Ab 2021 (Raumbetrachtung): Gutachten im Auftrag des Bundesministeriums für Wirtschaft Und Energie (BMWi) [Analysis of Regional Indicators For the Preparation of the GRW Funding Area from 2021 (spatial analysis): Expert opinion On Behalf of the Federal Ministry For Economic Affairs and Energy (BMWi)], in ifo Dresden Studien, Dresden, 2019, p. 83, ifo Institut2019.

[43] S. Maretzke, Infrastrukturindikator 2012: Ein Wichtiger Indikator für die Neuabgrenzung der Fördergebiete in Deutschland [Infrastructure Indicator 2012: an Important Indicator For the New Demarcation of the Funding Areas in Germany], in BBSR-Analysen Kompakt, 2014 no. 2014,5. Bonn: BBSR.

[44] Bundesministerium für Digitales und Verkehr, "Regionalstatistische raumtypologie (RegioStaR) des BMVI für die mobilitäts- und verkehrsforschung [Regional statistical area typology (RegioStaR) of the BMVI for mobility and transport research]." 2018. Accessed: Aug. 25, 2022. [Online]. Available: https://www.bmvi.de/SharedDocs/DE/Anlage/G/regiostar-arbeitspapier.pdf?_blob=publicationFile.

[45] T. Loga, N. Diefenbach, B. Stein, R. Born, Tabula Scientific Report Germany, IWU, 2012.

[46] Statistisches Bundesamt, "Bestand an wohnungen und wohngebäuden - Bauabgang von wohnungen und wohngebäuden - Lange Reihen ab 1969 - 2021 [stock of dwellings and residential buildings - Demolition of dwellings and residential buildings - Time series since 1969 - 2021]," 2021, [Online]. Available: https://www.destatis.de/DE/Themen/Gesellschaft-Umwelt/Wohnen/Publikationen/Downloads-Wohnen/fortschreibung-wohnungsbestand-pdf-5312301.pdf?_blob=publicationFile.

[47] Statistische Ämter des Bundes und der Länder, Wohngebäude Nach Baujahr [Residential Buildings By Year of Construction], Statistische Ämter des Bundes und der Länder | Gemeinsames Statistikportal, 2024. Accessed: Apr. 08[Online]. Available: https://www.statistikportal.de/de/wohngebaeude-nach-baujahr.

[48] "xgboost 1.7.6 documentation - categorical data." Accessed: Jul. 30, 2023. [Online]. Available: https://xgboost.readthedocs.io/en/stable/tutorials/categorical.html.

[49] N.V. Chawla, K.W. Bowyer, L.O. Hall, W.P. Kegelmeyer, SMOTE: synthetic minority over-sampling technique, Jair 16 (2002) 321–357, https://doi.org/10.1613/jair.953. Jun.

[50] "sklearn.Utils.Class_weight.Compute_sample_weight, " scikit-learn. Accessed: Aug. 15, 2023. [Online]. Available: https://scikit-learn/stable/modules/generated/sklearn.utils.class_weight.compute_sample_weight.html.

[51] J. Tanha, Y. Abdi, N. Samadi, N. Razzaghi, M. Asadpour, Boosting methods for multi-class imbalanced data classification: an experimental review, J Big Data 7 (1) (2020), https://doi.org/10.1186/s40537-020-00349-y. Art. no. 1Dec.

[52] "DummyClassifier, " scikit-learn. Accessed: Jan. 20, 2025. [Online]. Available: https://scikit-learn/stable/modules/generated/sklearn.dummy.DummyClassifier.html.

[53] "Understand your dataset with XGBoost — xgboost 1.7.6 documentation." Accessed: Apr. 18, 2024. [Online]. Available: https://xgboost.readthedocs.io/en/release_1.7.0/R-package/discoverYourData.html.

[54] "Regionaldatenbank Deutschland. GENESIS-Tabelle: 31211-02-01-4 gebäude mit wohnraum nach anzahl der wohnungen - Stichtag 09.05.2011, regionale Tiefe: Kreise und krfr. Städte, Gebäude- und Wohnungszählung 2011 (Zensus) [Regional database for Germany. GENESIS table: 31211-02-01-4 Buildings with living space by number of dwellings – reference date 09/05/2011, regional depth: districts and towns, 2011 census of buildings and dwellings (census)]." Accessed: Apr. 08, 2024. [Online]. Available: https://www.regionalstatistik.de/genesis/online.

[55] Statistische Ämter des Bundes und der Länder, "Regionaldatenbank Deutschland. GENESIS-Tabelle: 31211-03-01-4: gebäude mit wohnraum nach baujahr - Stichtag 09.05.2011 regionale tiefe: kreise und krfr. städte, Gebäude- und Wohnungszählung 2011 (Zensus) [Regional database for Germany. genesis table: 31211-03-01-4: buildings with living space by year of construction – reference date 09/05/2011 regional depth: districts and district-free cities, 2011 census of buildings and dwellings (Zensus)]." Accessed: Oct. 06, 2023. [Online]. Available: https://www.regionalstatistik.de/genesis/online.

[56] Umweltbundesamt, "Energieverbrauch privater haushalte," umweltbundesamt. Accessed: Apr. 29, 2024. [Online]. Available: https://www.umweltbundesamt.de/daten/private-haushalte-konsum/wohnen/energieverbrauch-privater-haushalte.

[57] Statistisches Bundesamt, "Umweltökonomische gesamtrechnungen - private Haushalte und umwelt [Environmental economic accounts – private households and the environment]," 2022.

[58] Kadaster, "Datamodel - BAG [data model - BAG]." Accessed: Feb. 17, 2025. [Online]. Available: https://bag.basisregistraties.overheid.nl/datamodel.

[59] S. Netherlands, Where Can I find District and Neighbourhood Data? Statistics Netherlands, 2025. Accessed: Feb. 17[Online]. Available: https://www.cbs.nl/en-gb/faq/infoservice/where-can-i-find-district-and-neighbourhood-data.

[60] Land Steiermark, ALS Gebäudemaske Steiermark [ALS Building Mask Styria], Open Government Data, Land Steiermark, 2024. Accessed: Nov. 15[Online]. Available: https://data.steiermark.at/cms/beitrag/11822084/97108894/.

[61] Amt der Tiroler Landesregierung, "Gebäude tirol [Buildings tyrol]." Accessed: Nov. 15, 2024. [Online]. Available: https://data-tiris.opendata.arcgis.com/maps/a1d63bd7edb34a76aca475d41e9a8ed9/about.

[62] Statistik Austria, "Datenangebot - regionalstatistische Raster [Data supply – regional statistical grids]." Accessed: Feb. 17, 2025. [Online]. Available: https://www.statistik.at/atlas/reg-datenkatalog/.

[63] F. Biljecki, Y.S. Chow, K. Lee, Quality of crowdsourced geospatial building information: a global assessment of openstreetmap attributes, Build. Environ. 237 (2023) 110295, https://doi.org/10.1016/j.buildenv.2023.110295. Jun.

[64] S. Thomas, D. Schüwer, F. Vondung, O. Wagner, Heizen Ohne Öl Und Gas bis 2035 – Ein Sofortprogramm Für Erneuerbare Wärme und Effiziente Gebäude. Im Auftrag Von Greenpeace e.V. [Heating Without Oil and Gas By 2035 – an Immediate Action Program For Renewable Heat and Efficient Buildings, On behalf of Greenpeace e. V.], 2022. Greenpeace, 2022. Accessed: Mar. 10[Online]. Available: https://www.greenpeace.de/publikationen/heizen-oel-gas-2035.

[65] Öko-Institut Prognos, Studie Im Auftrag Von Agora Energiewende, Agora Verkehrswende Und Stiftung Klimaneutralität [Climate-neutral Germany. Study commissioned By Agora Energiewende, Agora Verkehrswende and the Climate Neutrality Foundation], Wuppertal-Institut, "Klimaneutrales Deutschland, Prognos; Öko-Institut; Wuppertal-Institut, 2020. Accessed: Apr. 14, 2022. [Online]. Available: https://www.agora-energiewende.de/veroeffentlichungen/klimaneutrales-deutschland/.

[66] Fraunhofer IWES/IBP, Wärmewende 2030. Schlüsseltechnologien zur Erreichung der mittel- und langfristigen Klimaschutzziele im Gebäudesektor. Studie im Auftrag von Agora Energiewende. [Heat Transition 2030. Key technologies For Achieving medium- and long-Term Climate Protection Targets in the Building sector. Study commissioned By Agora Energiewende.], Fraunhofer IWES/IBP, 2017. Accessed: Mar. 25, 2022. [Online]. Available: https://www.agora-energiewende.de/fileadmin/Projekte/2016/Sektoruebergreifende_EW/Waermewende-2030_WEB.pdf.

[67] A. Mastrucci, A. Marvuglia, U. Leopold, E. Benetto, Life cycle assessment of building stocks from urban to transnational scales: a review, Renewable Sustainable Energy Rev. 74 (2017) 316–332, https://doi.org/10.1016/j.rser.2017.02.060. Jul.

[68] L.G. Swan, V.I. Ugursal, Modeling of end-use energy consumption in the residential sector: a review of modeling techniques, Renewable Sustainable Energy Rev. 13 (8) (2009) 1819–1835, https://doi.org/10.1016/j.rser.2008.09.033. Oct.

[69] Forschungsstelle für Energiewirtschaft e. V., "Wärmepumpen-Ampel [Heat pump traffic light]." Accessed: Mar. 08, 2023. [Online]. Available: https://waermepumpen-ampel.ffe.de/.

[70] K. Dabrock, K. Knosala, N. Pflugradt, H. Beuth, L. Kotzur, D. Stolten, The Potential of Combined PV and Air Source Heat Pump Systems in German Residential Buildings, presented at the EuroSun, 2022.

[71] C. Su, J. Dalgren, B. Palm, High-resolution mapping of the clean heat sources for district heating in Stockholm city, Energy Convers. Manage. 235 (2021) 113983, https://doi.org/10.1016/j.enconman.2021.113983. May.

[72] J.M. Miocic, M. Krecher, Estimation of shallow geothermal potential to meet building heating demand on a regional scale, Renew Energy 185 (2022) 629–640, https://doi.org/10.1016/j.renene.2021.12.095. Feb.

[73] M. Hörner and J. Bischof, "Typologie Der Nichtwohngebäude in Deutschland – Methodik, Anwendung und Ausblick [Typology of non-residential buildings in Germany – methodology, application and outlook]".