

MapLUR: Exploring a New Paradigm for Estimating Air Pollution Using Deep Learning on Map Images

MICHAEL STEININGER, KONSTANTIN KOBS, Albin ZEHE, and
FLORIAN LAUTENSCHLAGER, University of Würzburg, Germany
MARTIN BECKER, Stanford University, USA
ANDREAS HOTHOTH, University of Würzburg, Germany

Land-use regression (LUR) models are important for the assessment of air pollution concentrations in areas without measurement stations. While many such models exist, they often use manually constructed features based on restricted, locally available data. Thus, they are typically hard to reproduce and challenging to adapt to areas beyond those they have been developed for.

In this article, we advocate a paradigm shift for LUR models: We propose the **Data-driven, Open, Global (DOG)** paradigm that entails models based on purely data-driven approaches using only openly and globally available data. Progress within this paradigm will alleviate the need for experts to adapt models to the local characteristics of the available data sources and thus facilitate the generalizability of air pollution models to new areas on a global scale.

To illustrate the feasibility of the DOG paradigm for LUR, we introduce a deep-learning model called MapLUR. It is based on a convolutional neural network architecture and is trained exclusively on globally and openly available map data without requiring manual feature engineering. We compare our model to state-of-the-art baselines like linear regression, random forests and multi-layer perceptrons using a large data set of modeled NO₂ concentrations in Central London. Our results show that MapLUR significantly outperforms these approaches even though they are provided with manually tailored features.

Furthermore, we illustrate that the automatic feature extraction inherent to models based on the DOG paradigm can learn features that are readily interpretable and closely resemble those commonly used in traditional LUR approaches.

CCS Concepts: • **Computing methodologies** → **Neural networks**; *Image representations*; Supervised learning by regression; • **Applied computing** → **Environmental sciences**;

Additional Key Words and Phrases: Land-use regression, air pollution, deep learning

This work has been partially funded by the DFG grant “p2Map: Learning Environmental Maps—Integrating Participatory Sensing and Human Perception”

Authors’ addresses: M. Steininger, K. Kobs, A. Zehe, F. Lautenschlager, and A. Hothoth, University of Würzburg, Institute of Computer Science, Chair of Computer Science X, Am Hubland, 97074 Würzburg, Germany; emails: {steininger, kobs, zehe, lautenschlager, hothoth}@informatik.uni-wuerzburg.de; M. Becker, Stanford University, Stanford Medicine, Nima Aghaeepour Lab, 300 Pasteur Drive, Grant S280, Stanford, CA 94305, USA; email: mgbckr@stanford.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2374-0353/2020/04-ART19 \$15.00

<https://doi.org/10.1145/3380973>

ACM Reference format:

Michael Steininger, Konstantin Kobs, Albin Zehe, Florian Lautenschlager, Martin Becker, and Andreas Hotho. 2020. MapLUR: Exploring a New Paradigm for Estimating Air Pollution Using Deep Learning on Map Images. *ACM Trans. Spatial Algorithms Syst.* 6, 3, Article 19 (April 2020), 24 pages. <https://doi.org/10.1145/3380973>

1 INTRODUCTION

Air pollution is known to have adverse effects on human health and the environment [11, 22]. Thus, especially in areas with high population counts, it is important to control local pollution concentrations. For this reason, monitoring stations are deployed in many cities, which measure pollution continuously to assess whether the pollution levels are still within acceptable/legal limits. However, since the number of stations in a city is usually very limited, there are many areas where no air quality data is available. To fill these gaps, land-use regression (LUR) models are often used to estimate pollution concentrations in areas without monitoring stations [7, 26, 54].

Problem Setting. In recent years, a wealth of different land-use regression models have been developed that have shown to provide promising pollution estimates. However, these models (i) are at least partially based on neither globally nor openly available data [3, 7, 27] and (ii) often rely on hand-crafted features.

Thus, due to the local nature of the features, (i) these models usually do not generalize easily to locations other than the one they were developed for. Additionally, due to the involved hand-crafting process, (ii) optimizing the features for new models in specific study areas is a cumbersome process.

Approach. To address the challenges inherent to inaccessible data and manual feature engineering in land-use regression models, in this work, we advocate a paradigm shift toward purely data-driven land-use regression models based on open and globally available data. We call the corresponding paradigm **DOG (Data-driven, Open, Global)**. More specifically, models adhering to DOG work directly on raw data, automatically extracting their features from the input. While such data-driven methods have proven successful in multiple application domains [33], they have so far not been introduced to land-use regression. Land-use regression models following the DOG paradigm have multiple advantages: (i) they can be fit more easily to different study areas than other, more specialized land-use regression approaches, (ii) they do not require manual feature engineering, and (iii) they can be reproduced by other researchers without requiring access to data sources that are not easily available. To demonstrate the feasibility of this paradigm, we introduce the MapLUR model. MapLUR implements DOG by using deep learning, specifically a convolutional neural network architecture. It automatically extracts features from map images, which are openly available almost anywhere on earth, and estimates air pollution based on these features.

Experimental Evaluation. We assess the performance of MapLUR by comparing it against state-of-the-art land-use regression models like linear regression, Random Forests (RF), and Multi-layer Perceptrons (MLP) on modeled NO₂ concentration data from the London Atmospheric Emissions Inventory (LAEI) [2]. In the process, we employ different types of images including map images from OpenStreetMap and Google Maps [24] as well as satellite imagery from Google Maps. We find that our model works best using map images from OpenStreetMap and that it outperforms all baselines significantly.

Furthermore, we analyze the data requirements of MapLUR and the baselines. We find that common for deep-learning models, MapLUR requires more training data than models that rely on

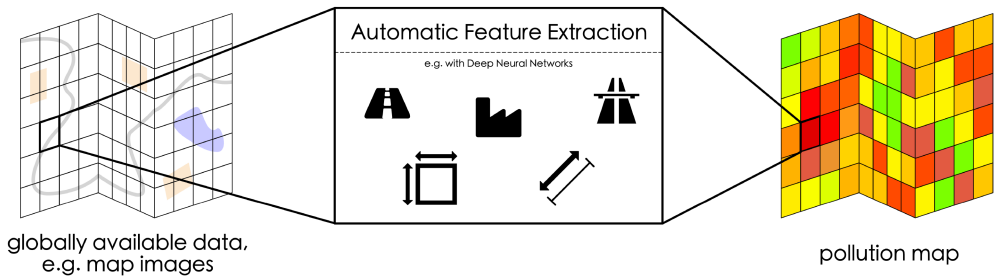


Fig. 1. MapLUR: Automatic feature extraction and globally available data for air pollution modeling. We propose the DOG paradigm (Data-driven, Open, Global) for land-use regression, which advocates openly and globally available data, and automatically extracting features to estimate pollution. Following this newly introduced paradigm, we propose the MapLUR model. It consists of a deep neural network architecture that estimates pollution concentrations for specific locations directly from globally available map images (rendered maps or satellite images), resulting in area spanning pollution maps. For this, MapLUR automatically learns to extract features from the given map images. The extracted features closely resemble manually engineered features for land-use regression models.

hand-crafted features. In this context, we evaluate how far the training set can be reduced and discuss possible approaches to further address this challenge.

Finally, we analyze the automatically extracted features by observing which parts of the map images were particularly important for our model using guided backpropagation [51] and artificial map images. The analysis shows that the learned MapLUR features strongly relate to hand-crafted features as commonly used in land-use regression models.

Contribution. Our core contributions in this work are threefold:

- (1) We propose DOG, a new, data-driven paradigm to land-use regression. Models following this paradigm should not require manual feature engineering and only rely on openly and globally available data sources.
- (2) We introduce MapLUR, a land-use-regression model based on DOG. MapLUR employs a deep-learning approach to automatically extract features from map images. We show that this model is able to outperform traditional land-use regression models when trained on a sufficiently large data set.
- (3) We demonstrate that, contrary to popular believe, models based on the data-driven paradigm are not necessarily black-boxes by inspecting the features MapLUR extracts, finding that the automatically extracted features strongly relate to typical manually engineered features for land-use regression models.

Structure. This work is organized as follows. Related work is summarized in Section 2. Section 3 describes the air pollution data and image data used in this work. DOG and the MapLUR model are introduced in Section 4. The experiments and the baseline models are described in Section 5. Section 6 presents our results and analyzes our model. We discuss advantages and limitations of DOG and MapLUR in Section 7. Finally, Section 8 concludes this work.

2 RELATED WORK

Land-use regression has been an active field of research for many years now. Work done in the previous decade has laid important foundations for current land-use regression models and established linear regression techniques as the de facto standard model [7, 43, 54]. Especially noteworthy is the Escape project [7, 20], which built models for 36 European areas. The model building

procedure of this project has become a standard approach [40, 41, 43, 52, 53]. To make the application of land-use regression models easier, there is a tool available that automizes the process of variable generation, modeling and prediction with a model based on linear regression [42].

However, more advanced machine learning methods are starting to become more common. One example for these approaches are Random Forests [9]. They have been used successfully to estimate elemental components of particulate matter in Cincinnati, Ohio [10] and NO₂ pollution in Geneva [14].

Another example for a more advanced method are neural networks. These models have been used to estimate a range of pollutants successfully, as shown in various publications. For example, they have been applied to NO₂ [1, 14, 37], PM_{2.5} [1, 5, 6, 23, 55], PM₁₀ [3, 37, 55], and surface dust [12] concentrations. The neural-network-based models are typically simple multi-layer Perceptrons. However, there are deep-learning models that use recurrent neural networks or deep belief regression networks. These models differ from this work in that they are used to forecast pollution concentrations from earlier measurements or fill missing values for locations where measurements already exist [5, 6, 23, 55], while we estimate pollution for locations without measurements. To the best of our knowledge, there are no deep-learning models for our setting. Both Random Forests and neural networks have been shown to outperform linear regression in land-use regression [10, 14].

Support vector regression models [19] are another possible approach. There are models that can forecast pollution concentrations using this technique [35, 47], but there do not seem to be any land-use regression models with this type of model.

All aforementioned land-use regression models rely on manually engineered features, which are typically gathered from various locally available data sources that might not be available elsewhere. In contrast to all methods above, we propose a deep-learning model based on convolutional neural networks (CNNs), which is able to automatically learn relevant features from openly available maps.

Such image-based approaches have been used before in the context of air quality estimation and pollution detection. Singh [50] interpreted modeled air pollution data as images and used non-machine-learning image classification techniques to detect higher pollution episodes. Furthermore, CNNs have been used before in the context of air quality estimation by Zhang et al. [57] and Li et al. [36], who proposed models to estimate air haze level using photos from, for example, mobile phones or webcams. In contrast, our work uses map and satellite imagery depicting land-use as model input, making our model more closely related to land-use regression models. Additionally, our model estimates pollution concentrations instead of haze levels.

3 MATERIALS

In this section, we introduce the air pollution data set we use to train and evaluate our method as well as the data sources from which we extract map and satellite images.

3.1 Air Pollution Data

We train and test our model using pollutant concentrations from the London Atmospheric Emissions Inventory (LAEI) [2]. It contains modeled annual mean concentrations of NO₂ and PM₁₀, among other pollutants, at a 20 m grid level for the complete Greater London area in 2013. For our main model development and evaluation, we use the NO₂ concentrations of the data set, since it is a very frequently used pollutant for land-use regression models. The data is the result of a dispersion model that incorporates a vast number of input factors like for example road and rail networks, traffic data, aviation, pollution from individual industrial premises, domestic and commercial fuel

Table 1. Statistics of London's NO₂ Concentrations

Data set	Count	Mean	SD	Min	Max
Central London	113,680	50.90	15.02	37.12	253.89
Sampled subset	4,500	50.85	15.02	37.17	171.06

The statistics for Central London include all LAEI cells in the bounding box. The sampled data set contains randomly sampled cells from the Central London data set, which are used for training and testing models in our experiments. Mean, standard deviation (SD), minimum values (Min), and maximum values (Max) are in $\mu\text{g}/\text{m}^3$.

consumption, as well as fires. Through this approach, 5,856,428 data points were generated where each data point represents a 20 m by 20 m cell [2].

We sample a training set consisting of 3,000 data points and a test set consisting of 1,500 data points from the Central London part of the data set to have a reasonable number of urban data points for our experiments. We choose data points from Central London, because we believe that it is more important in practice to reliably estimate pollutant concentrations in highly polluted areas with a large population than in more rural areas. For this, we define a geographical rectangle that roughly describes Central London and only use cells within. The box's north western corner is at 526,660 easting and 183,220 northing while the south eastern corner is at 534,760 easting and 177,640 northing specified in British National Grid coordinates. The sampled data points are depicted in Figure 8 in Appendix A. Descriptive statistics for the cells in Central London and the sampled subset used for training and testing can be found in Table 1.

The map images, which we use to depict the areas of the data points, show 80 m by 80 m areas even though the air pollution data is available at a 20 m grid level. The 20 m by 20 m cells are in the center of these images. This allows MapLUR to see more of the surroundings and incorporate information about distant emission sources. To avoid a potential evaluation issue, we sample data points in such a way that no images can overlap. Any overlap could lead to a situation where the model already roughly knows the pollution concentration for a test data point, since it might recognize the test data point's area from the image of a nearby training data point. Such implicitly learned proximity of data points could give our model an unfair advantage, which we avoid with our procedure.

3.2 Image Data

This section describes the sources of map and satellite image data that we use in this article as well as the preprocessing applied to the images to generate training samples.

3.2.1 Image Sources. There is a variety of globally available sources for map images, two popular services being Google Maps [24] and OpenStreetMap [46]. While Google Maps is a commercial and proprietary service, OpenStreetMap is an open database for map data that is built and maintained by volunteers. Data from OpenStreetMap can be used to render maps in various ways through different stylesheets. In this work, we render map images based on OpenStreetMap data using a slightly modified version of the default stylesheets used on the official OpenStreetMap website. It differs from the default in that we do not render text like street or station names, since labels obstruct map features making them harder to recognize and often only carry very localized information, thus, possibly reducing generalizability. For tile rendering, we use *mod_tile*, which is a module for the Apache web server with the rendering back-end *renderd* [45]. In addition to OpenStreetMap images, we use map and satellite images from Google Maps, to compare the effectiveness of each visualization for this task. Since Google Maps is proprietary, the images cannot be

easily customized to the same extent as OpenStreetMap images. We therefore use them without modification.

3.2.2 Image Preparation. Before using MapLUR it is necessary to prepare map or satellite images. We found through preliminary experiments that images depicting 80 m by 80 m provide the best performance in this setting, as can be seen in Appendix B. To depict the correct area in an image, we approximate the 80 m distances using a meter per pixel value that depends on the zoom level of the image and the latitude of the location due to the Mercator projection, which both OpenStreetMap and Google Maps use. We obtain the images at zoom level 17, resulting in a pixel extent of approximately 0.75 m by 0.75 m at London's latitude. Thus, the images have a resolution of 106 px by 106 px. The rendered images are then scaled to a fixed resolution of 224 px by 224 px similar to what popular CNN architectures for the ImageNet competition [18] use. This way, we avoid having to change the model architecture when depicting different sized areas in the images or when we fit the model to locations with a different latitude, which would also result in differently sized images. Using a model input resolution that fits the images exactly could reduce the model size and improve training and inference speed. Nonetheless, we found it unnecessary considering the already acceptable speed and we favored the increased flexibility.

4 METHODS

In the following, we introduce our data-driven paradigm DOG to land-use regression, which suggests that models should estimate air pollution by automatically extracting relevant features from openly and globally available data. We also present our model *MapLUR*, which follows this paradigm by using a convolutional neural network as an automatic feature extractor, taking as input globally and openly available map images.

4.1 The DOG Paradigm

Our data-driven paradigm DOG (Data-driven, Open, Global) aims to alleviate the issues of manual feature engineering as well as only locally applicable models. To this end, it requires models to fulfill the following criteria:

- **Automatic extraction of relevant features:** Models should learn a function from the raw input to the desired output, which leads to automatic development of features and may even uncover relevant factors that are not yet known as having an influence on air pollution. This can for example be achieved by deep-learning methods like the convolutional neural network introduced in the next section.
- **Usage of globally available data sources:** Models should rely exclusively on data sources that are available for (almost) all parts of the world. This allows the ubiquitous application of the model without collecting additional, only locally available data sources.
- **Usage of openly available data sources:** Models should rely exclusively on openly available data sources. This allows researchers to reproduce and improve the model's results without requiring access to paid or not publicly available resources.

4.2 The MapLUR Model

In this section, we propose the specific model MapLUR based on the DOG paradigm described above. Our model applies the paradigm by using map images as a globally, openly available source of information and extracting features through deep learning, more specifically a convolutional neural network (CNN). This type of network is a natural fit for our setting, since it is specifically designed to work with two-dimensional shapes like images. CNNs utilize spatial locality within the images and reduce the number of learned parameters through weight sharing.

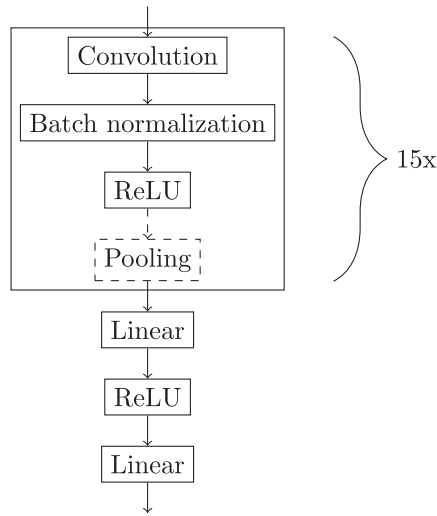


Fig. 2. Structure of MapLUR. The model consists of 15 feature-learning building blocks, which contain a convolution layer, batch normalization, rectified linear units (ReLU), and sometimes a pooling layer. These building blocks are concatenated and only the 1st, 3rd, 5th, 7th, 10th, and 13th block contain a pooling layer. These blocks are followed by a simple fully connected layer with ReLU activation and finally a single fully connected neuron with linear activation that returns the estimation of the pollution at the given location.

These concepts make it feasible to learn relevant features from raw images, in contrast to fully-connected networks, which would need an impractical amount of parameters [34].

The structure of MapLUR is depicted in Figure 2. It contains 15 convolutional layers with batch normalization [29] and rectified linear units (ReLU), which are pair-wise linear activations [44]. The last convolutional layer is followed by a fully connected layer with 128 neurons and ReLU activation (depicted as the third and second to last boxes in Figure 2). These neurons are then connected to a single neuron with linear activation that produces the estimated pollution concentration. Each convolutional layer has 16 filters, a kernel size of 3, a padding of 1, a dilation of 1, and a stride of 1. The output size of these layers is the same as their input size. Maximum pooling layers with a kernel size of 2 and stride of 2 are applied after the ReLUs of the 1st, 3rd, 5th, 7th, 10th, and 13th convolutional layer to reduce the number of activations. We found this architecture and the corresponding hyperparameters by evaluating different variations of the model using ten fold cross-validations on the training set.

We use ReLU activations, since they have shown to work well for many different tasks and models, making them the most popular activation function for deep-learning applications [33]. While trying different architectures, we have also experimented with SELU [31] and RReLU [56] activations, but we found no improvements with these functions. The linear activation in the final layer is common practice for regression tasks [32]. It does not restrict the range of resulting values allowing the model to estimate any value.

5 EVALUATION

To evaluate MapLUR, we conduct several experiments and compare our model to baseline models that are commonly used in land-use regression. The experiments and the baseline models are described in the following.

5.1 Experimental Setting

We conduct four experiments using MapLUR, varying the data available to the model. For all experiments, MapLUR is trained using the Adam optimizer [30] on batches of size 400 for at most 2,000 epochs with a learning rate of 0.0001. We augment the training data by flipping or transposing the images. Additionally, we employ early stopping, interrupting training when the validation performance has not increased for 20 epochs in a row.

Experiment 1—OpenStreetMap. In the first experiment, only OpenStreetMap images are used as input to the CNN. The input images have three channels (RGB), 224 px by 224 px, and depict an area of 80 m by 80 m. All labels were removed from the rendering process of the images, as described in Section 3.2.1.

Experiment 2—Google Maps. Instead of OpenStreetMap images, Google Maps images were captured and fed into the model in this experiment. The same size as in the previous experiment was used. As Google Maps data is proprietary, modifications cannot be made as easily and to the same extent as with OpenStreetMap images. Therefore, text labels are present in the imagery.

Experiment 3—Google Maps Satellite. For the third experiment, instead of stylized map images, we use satellite images from Google Maps Satellite [24], which uses imagery from both satellites and aerial surveys. Training and test images from Google Maps Satellite have the same size and zoom levels as the OpenStreetMap and the Google Maps images.

Experiment 4—OpenStreetMap and Google Maps Satellite. Experiment 4 then combines OpenStreetMap images and satellite images by concatenating the two RGB images to one six-channel tensor.

5.1.1 Evaluation Setup. The models are evaluated using standard metrics for the evaluation of land-use regression models, namely, R^2 and RMSE. Both metrics are explained in Appendix C. In all experiments, the model is initialized and trained 40 times on the training set and evaluated on the test set, both of which are described in Section 3.1. The average of the resulting evaluation metrics is then used as the final score to counteract unfortunate initialization results. Additionally, the sample of 40 evaluation runs can be used as the input to statistical significance tests to formally confirm differences in evaluation results.

5.2 Baselines

To determine how well our model works, it is necessary to compare it to other methods. Therefore, we first describe a set of features that is used by our baseline models. Thereafter, four baseline models are introduced, namely, a mean baseline, linear regression, Random Forest, and multi-layer Perceptron. The last three of the aforementioned models are commonly used in land-use regression. Random Forests and multi-layer Perceptrons tend to yield state-of-the-art results as described in Section 2.

5.2.1 Features. We use a set of standard land-use and road-related features for our baseline models. These features have shown to be important influencing factors for air pollution [20]. All of these features can be calculated from OpenStreetMap data, since we want to provide similar information to all models for a fair comparison. The features include the areas of commercial, industrial, and residential land-use, the lengths of big and local streets, and the distances to the next traffic signal, motorway, primary road, and industrial premise. Big streets include streets that are classified as either motorway, trunk road, primary road or secondary road in OpenStreetMap while all other streets are local streets. Most features are typically calculated for different *buffers*, which

are areas with a specific radius around data points. For example, the areas of different types of land-use and the lengths of streets are calculated for 50 m and 100 m buffers to give the baseline models similar sight into the surroundings as the MapLUR model. However, the features that calculate the distance from each data point to specific locations like the closest traffic signals, roads or industrial premises exist only once and are not calculated for different buffers. Due to these features, the baseline models are given a slight advantage, since they can get information from entities that are further away than MapLUR can see.

5.2.2 Mean. A simple baseline for a regression task is the mean baseline. It disregards all features and estimates the mean value of all training data points for each test data point. This baseline provides performance values that every other model should beat.

5.2.3 Linear Regression. The most common approach to land-use regression is linear regression. Therefore, it is useful to compare our novel model to this type of model.

We use the same supervised stepwise selection as Eeftens et al. [20] for selecting the most relevant subset of features. A description of this procedure can be found in Appendix D.1. After applying the stepwise selection on the development set the model is left with the variables *length of big streets* (50 m buffer), *distance to the next industrial premise*, and *distance to the next traffic signal*.

5.2.4 Random Forest. The Random Forest is a more powerful model that was shown to work well for land-use regression and can often provide better performance than typical linear regression approaches, as it can model non-linear correlations between features [10, 14]. Therefore, we use it as another baseline in this work.

This model is built in a similar way to the procedure in Brokamp et al. [10]. Details are in Appendix D.2. The final Random Forest model uses the variables *distance to the next industrial premise*, *distance to the next primary road*, *distance to the next traffic signal*, *distance to the next motorway*, *length of big streets* (50 m buffer), and *area of residential land-use* (100 m buffer). It builds 394 trees using bootstrap samples, considers at most 42.79 % of the available features per split, needs at least three samples to split a node, and needs at least three samples for a leaf node.

5.2.5 Multi-layer Perceptron. Neural networks, or more specifically multi-layer Perceptrons (MLPs), are models whose popularity for land-use regression tasks has grown in recent years and that often outperform other baselines [1, 3, 12, 14, 37]. Additionally, evaluating multi-layer Perceptrons (even if not directly applied to image data) illustrates the performance of neural networks that, in contrast to the MapLUR model, are not based on convolutions.

Again, we follow the model development procedure of previously published work. In this case, we base our procedure on the one used by Alam and McNabola [3]. Appendix D.3 contains a description of this procedure.

The MLPs use the variables *distance to the next industrial premise*, *distance to the next primary road*, *distance to the next traffic signal*, *distance to the next motorway*, *length of big streets* (50 m buffer), *length of local streets* (100 m buffer), *area of industrial land-use* (100 m buffer), *area of commercial land-use* (100 m buffer), and *area of residential land-use* (100 m buffer). The architecture search found that the best performing MLP model has a single hidden layer with 29 neurons. This is similar to the MLPs used by previous publications [1, 3, 12, 14, 37].

6 RESULTS AND ANALYSIS

Given the baseline methods and MapLUR's description, we now present the results for our experiments and analyze MapLUR in terms of data requirements and features learned.

Table 2. Results of Baseline Methods and Experiments

Model	R ²	RMSE [$\mu\text{g}/\text{m}^3$]
Mean baseline	0.000	13.971
Linear regression	0.487	10.004
Multi-layer Perceptron	0.499	9.887
Random Forest	0.662	8.119
MapLUR experiment 1: OpenStreetMap	0.673	8.002
MapLUR experiment 2: Google Maps	0.537	8.918
MapLUR experiment 3: Google Maps Satellite	0.206	12.389
MapLUR experiment 4: OpenStreetMap and Google Maps Satellite	0.660	8.112

MapLUR with OpenStreetMap images is providing the best performance overall, beating all baselines and all other MapLUR variants. Using satellite images from Google Maps instead of OpenStreetMap images decreases the metric scores on the evaluation set. Combining both image types does not improve the score. The Random Forest model is outperforming all other baseline models on this data set, which makes it the best baseline. All results are significantly different to each other.

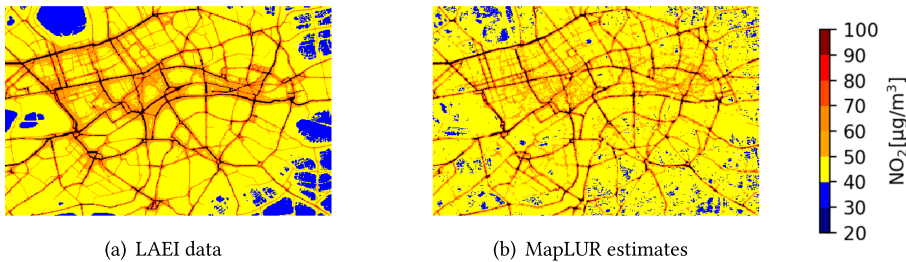


Fig. 3. Comparison of LAEI data and MapLUR estimates. The estimated map shares strong similarities with the original data. It can be seen that MapLUR is able to recognize streets and accurately associate them with high pollution concentrations. However, our model tends to overestimate pollution in areas with very low pollution concentrations and underestimate pollution for areas with no road close by.

6.1 Experiments

Table 2 shows the results of the baseline methods as well as MapLUR's results for our experiments. All results in the Table are significantly different from each other. To verify this, the metrics of each model are tested for normality using the test from D'Agostino and Pearson [15, 16] with $p < 0.05$. The statistical significance for models whose metrics are normally distributed are tested using a t-test, while the other models' metrics are tested with the Wilcoxon signed-rank test, both testing for $p < 0.05$. Additionally, Bonferroni correction [8] is applied, which further substantiates the statistical significance, since $p < \frac{0.05}{n}$ with $n = 7$ for each model pair. $n = 7$ is chosen to account for the number of hypotheses that are tested on the same data (each model is tested against 7 other models).

Baselines. As described before, we use a simple mean baseline, a linear regression, an approach with Random Forests, and a multi-layer Perceptron with manually engineered features from OpenStreetMap data. The results in Table 2 show that the Random Forest is performing considerably and significantly better than both the linear regression and the MLP.

Experiment 1—OpenStreetMap. Our model with OpenStreetMap images performs better than all baselines regardless of metric, which can be seen in Table 2.

Figure 3 shows the original NO₂ concentrations of the LAEI data set in Central London and the estimates of MapLUR. It can be seen that our model is able to come rather close to the original data using only OpenStreetMap images, but tends to overestimate values of areas with low pollution concentrations and underestimate values of areas that are not in the vicinity of roads.

Experiment 2—Google Maps. This experiment uses Google Maps imagery instead of OpenStreetMap images. A drop in R² of more than 10 percentage points in comparison to the previous experiment and a higher RMSE value may be explained by the styling of Google Maps images. Google Maps contain fewer color-coded entities. Especially streets, that are a common entity for land-use regression features, are not diversified as much as in OpenStreetMap images. The differences can be seen in Appendix E.

Experiment 3—Google Maps Satellite. This experiment uses Google Maps Satellite imagery as input. Table 2 shows that using only Google Maps Satellite imagery leads to a considerable drop in performance, even worse than the linear regression baseline with an R² of 0.206 and an RMSE of 12.389.

These results are most likely due to the noise in the satellite images, which makes it harder to discern influencing factors for air pollution. The hand-labeled map images therefore help a lot as they already encode the desired entity labels as colors.

Experiment 4—OpenStreetMap and Google Maps Satellite. The last experiment combines map and satellite imagery by concatenating both three-channel RGB images to one six-channel tensor. OpenStreetMap images are used for the map images, since they have shown better performance than Google Maps in our task. Satellite imagery is taken from Google Maps Satellite. As both image sources use the same spatial resolution of 80 m by 80 m, local OpenStreetMap data should only be augmented by the satellite images. However, the results in Table 2 show no performance gain compared to using OpenStreetMap only. In fact, the results are worse and significantly different for both models.

Computation times. The computation times for the various models evaluated vary due to the different model complexities. For example, with an Intel Xeon E5-2690V4 CPU training the MLP takes on average 25 seconds, Random Forest trains on average for 11 seconds, and the linear regression model is typically built in a single second. Estimating pollution concentrations with these simple trained models for an area like Central London only takes seconds. MapLUR is the most complex model among them, but can be trained in reasonable time on commodity hardware. We found that training it on a single consumer graphics card (Nvidia GTX 1080 TI) takes about 50 minutes. Once MapLUR is trained it can estimate pollution concentrations for the complete Central London area in 35 seconds.

6.2 Model Analysis

After seeing that MapLUR can work well, we now further analyze the model. First, we assess the data requirements of MapLUR. Then, we demonstrate that our model can be made interpretable by analyzing what our model has learned through guided backpropagation and by creating fake OpenStreetMap images.

6.2.1 Analyzing Data Requirements. The previous experiments showed that our model can successfully model air pollution in a data-driven way. While this is the main focus of this article, the data we use (3,000 data points) is larger than those typically available in a real world setting. In the following, we analyze the actual data requirements of all models evaluated above. Overall, the

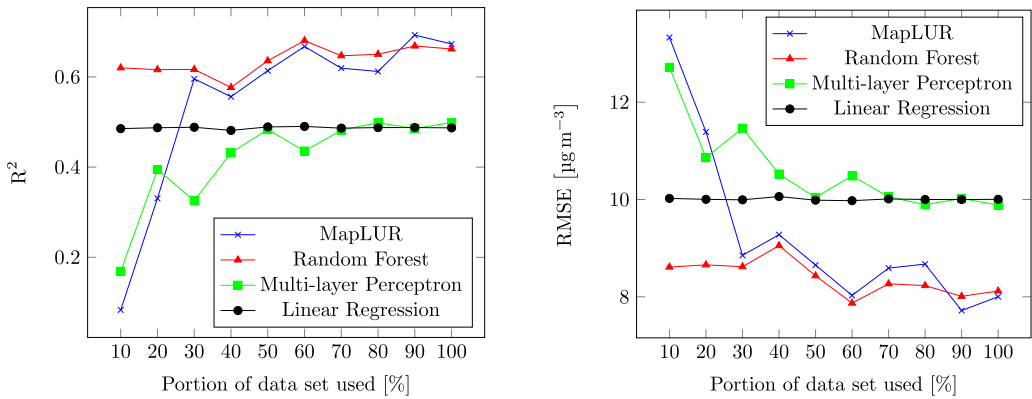


Fig. 4. Data requirements of land-use regression models. These graphs show the performance of commonly used land-use regression models and MapLUR with varying training data set sizes. The portion of the data set used refers to the size of the NO_2 data set with which the models were trained in our experiments. Thus, 100% is equivalent to 3,000 data points. Each point in the graphs is the mean of 40 model runs, except for MapLUR's points, which we only ran five times per data set size due to the model's computational complexity. This shows that MapLUR can provide comparable results with 900 data points and it tends to improve with more data. Multi-layer Perceptrons behave similarly but they need more data to reach other baselines. Linear Regression and Random Forests are less dependent on data set size.

corresponding results will inform future studies on data requirements and point toward necessary methodological advancements.

When gradually reducing the number of training data points, we noticed a drop in performance with smaller data sets for all models except for linear regression (cf. Figure 4). Models based on neural networks experience a more pronounced performance loss in comparison to, for example, Random Forests, where there is only a slight decline. We believe that this stems from the size and complexity of these models compared to simpler models like linear regression models. More parameters need to be trained, which tends to require more training examples. However, about 30 % of the training data is still sufficient for our model to exhibit performance comparable to the strongest baseline, which is the Random Forest trained on hand-crafted features.

Addressing the increased need for data is an important point for future work, which we also discuss in some more detail in Section 7.

6.2.2 Understanding Estimates using Guided Backpropagation. In this section, we apply a technique called *guided backpropagation* [51], which allows us to visualize the regions of the input image that the network focuses on for its estimation. This approach starts off with a forward pass of an image. Thereafter, the gradient of the activation is computed with respect to the input image. At each ReLU in the model, positive gradients whose corresponding output during the forward pass was negative and negative gradients are set to 0 so that only features that contribute to the estimated pollution concentration are shown. This allows us to visualize which parts of the image the model is paying attention to. Several examples are shown in Figure 5.

The guided backpropagation shows that the model is paying special attention to motorways, trunk roads, and primary roads, which are rendered in red or orange colors in OpenStreetMap. This shows that MapLUR is able to automatically learn intuitively relevant features, since traffic is known as a large factor for NO_2 pollution [13]. MapLUR also considers buildings, foot paths and cycle paths to some extent for its estimates while it seems to be ignoring water and park areas. The model tends to pay more attention to pixels close to the center, which is understandable, since

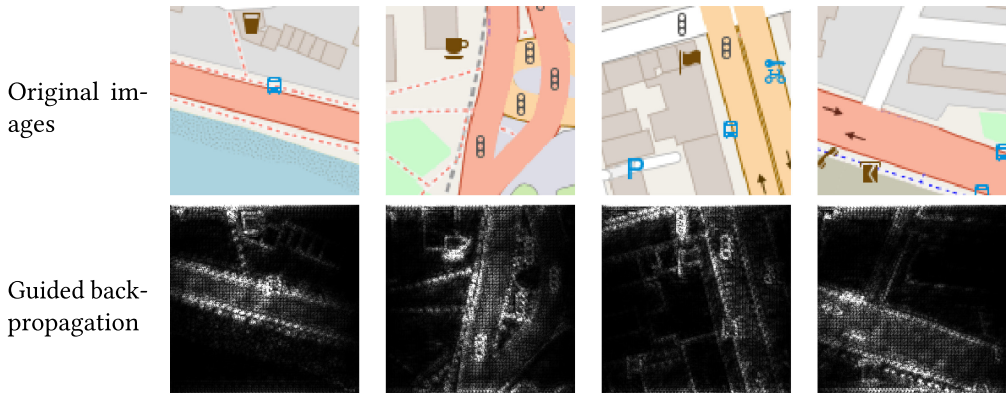


Fig. 5. Visualization of detected features using guided backpropagation [51]. This technique highlights important pixels in the input images by visualizing gradients of the activation with respect to pixel intensities. All negative gradients and positive gradients, whose corresponding output during the forward pass was negative, are set to 0 at each ReLU during backpropagation. This approach reveals parts of an image that contribute to the pollution. As one would expect, the model is concentrating on large streets. (Original images: ©OpenStreetMap contributors)

we estimate the pollution concentration for the 20 m by 20 m areas that are in the center of each image.

6.2.3 Analyzing Entity Influence using Artificial Map Tiles. One of the biggest advantages of using a DOG-based model for air pollution estimation is that it extracts features by itself, while previous work always used hand-engineered features. The leading question in developing land-use regression methods in previous work is: *What entity of what area in what distance to the center is contributing to the pollution?* From this, three categories of features arise: entity features, area features, and distance features. We now want to investigate the correlation of these features with the model's output. For this, we take advantage of the well-defined structure of map images with different color-coded entities. Map images therefore can easily be recreated using graphic editing software, which makes it possible to create artificial OpenStreetMap images for which we can control the features separately while keeping all other features fixed. We then observe changes to the model's output while modifying the values of these features.

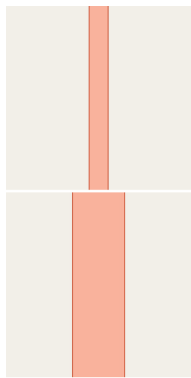
Entity Features: Entity features describe *what* is seen on the image. Entity features are often used for the estimation of air pollution, as, for example, industrial areas are usually contributing more to air pollution values than parks. In this experiment, we investigate how certain entities are influencing the model estimate. We build two kinds of images: On the one hand, we create images that are each completely covered by one specific type of entity, resulting in uniformly colored square images. On the other hand, the same images are then overlaid by the depiction of a motorway and a trunk road. We expect that different underlying entities provide different estimates according to the usual presence of sources for NO_2 pollution. We also expect an increase in the air pollution estimate whenever a road is added to the underlying entity. Depending on the type of road this increase might fluctuate. Table 3 shows the resulting pollution estimates by the CNN.

Different underlying entities do not lead to large differences in pollution estimates if there is no road. Only completely covering the image by a motorway or trunk road results in an estimate of over $45 \mu\text{g}/\text{m}^3$. Additionally, trunk roads seem to have a much higher impact on the air pollution

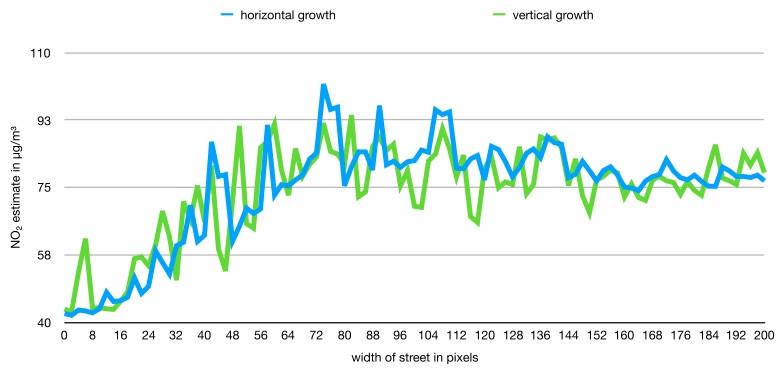
Table 3. Model Estimate for a Given OpenStreetMap Entity

Road type	Entity Name								
	Industrial area	Residential area	Commercial area	Park	Forest	Water	Neutral	Motorway	Trunk
No road	37.71	38.29	38.72	38.87	39.27	41.66	42.06	47.23	80.63
Trunk	61.70	50.94	59.73	64.14	57.94	59.62	58.62	—	—
Motorway	60.04	48.48	64.18	46.05	54.21	53.45	55.00	—	—

The entities span across the whole image and they are overlaid with different types of roads. Overlaying a road with another road does not make sense so these values are omitted. The “neutral” entity is a background that is used by OpenStreetMap for indicating land with no particular land-use. All estimates are in $\mu\text{g}/\text{m}^3$.



(a) Two fake examples with a trunk road of width 10 px and 60 px, respectively.



(b) NO_2 estimate for a given width of the road.

Fig. 6. Varying the width/area of the street while keeping other features such as distance to the center and type of street fixed.

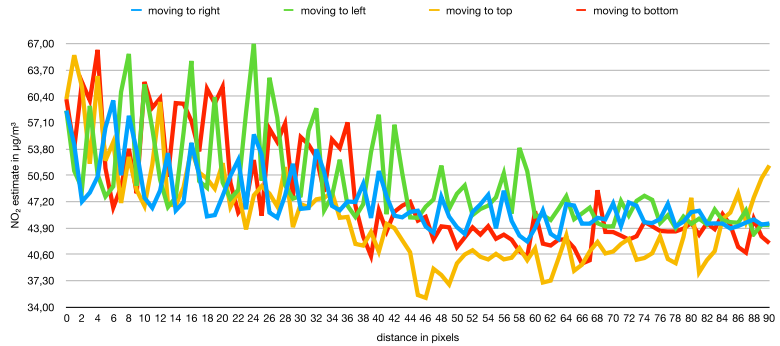
estimate than motorways. Adding a trunk road or motorway to any entity increases the air pollution estimate as expected. The amount of increase depends on the underlying entity of the map and what kind of road is present. This shows that the relationship of the entities that are visible in the map image are also important. There seem to be complex correlations between different entity features, which cannot be modeled easily in simpler models like linear regression.

Area Features: Area features describe *how large* a given entity is in the image. The area that an OpenStreetMap entity has on an image should contribute to the estimated pollution value. In this experiment, we use the trunk road entity to show the influence of the area. We build multiple images that contain a straight road that goes top to bottom or left to right through the center of the image. As the background, we always use the same neutral background that depicts general land-use in OpenStreetMap. We then vary the width of that street either horizontally or vertically, depending on the street direction. A linear increase in the street’s width is equivalent to a linear increase in the street’s area.

Figure 6 shows some of the artificial OpenStreetMap images as well as a plot of MapLUR’s output given the street width in pixels. As expected, an increasing width—and therefore an increasing area—of the street tends to increase the pollution estimate. Both horizontal and vertical growth have very similar curves that are not linear but instead seem to be more logarithmic. The similarity



(a) Two fake examples with a trunk road that is moved by 10 px and 60 px, respectively.



(b) NO_2 estimate for a given distance from the center of the image to the trunk road.

Fig. 7. Varying the distance of a trunk road to the image center pixel by pixel while keeping other features fixed.

was expected, as during training, the images are augmented by rotation and flipping such that the direction of streets should not have any impact on the overall output.

Distance Features. Distance features describe *how far away* a given entity is from the image center. For this experiment, we create images that contain only one straight trunk road that is then moved vertically or horizontally, depending on the direction of the street. With this setup, we can control the distance of the motorway to the center of the image while fixing the area and entity features. We expect that the model produces higher estimates for images where the street is closer to the center, as this behavior was already observed in the guided backpropagation results. Also, the desired value from LAEI is coming from a 20 m subframe of the image, which is in the image center. The model therefore should have learned a tendency to weight features from the center of the image more than from the borders.

Figure 7 shows image samples and the resulting estimation curves when moving the trunk road farther away from the center of the image. As expected, the proximity of a street to the image's center contributes to the overall NO_2 estimate positively. Pearson correlations of the distance with the estimated values are always lower than -0.6 , indicating a relatively strong negative correlation. The curves that are shown are also not linear and can be better fitted by polynomials with a squared feature term than by a line. To capture this non-linearity, more sophisticated methods need to be used, which justifies the use of Random Forests or neural networks.

7 DISCUSSION

In this section, we discuss some advantages and some current limitations of our proposed paradigm and model. We believe that the advantages provide valuable additions to the models currently applied in land-use regression. Since this is the first article applying a model based on our new paradigm, there are still some limitations regarding the applicability of our model to real world data sets, for which we provide some possible ways to overcome.

Interpretability. First, purely data-driven models tend to be harder to interpret than simpler models, which is why they are often thought of as black-box models. This also raises the concern

that the models may put too much focus on unreliable features that explain the specific data set well, but fail to generalize to other data sets. To alleviate these concerns, we have shown that it is possible to reveal the inner workings of the model, finding that the model's output heavily relies on land-use features such as streets or commercially used areas. These features are also commonly employed by traditional land-use regression models. We have also shown that MapLUR implicitly focuses on other commonly used land-use information such as distance and area features. This illustrates that purely data-driven approaches can yield interpretable models.

Feature extraction and complexity. In addition to information that closely resemble hand-crafted features, our model is able to extract signals from the image in an automated and optimized fashion. Thus, it can potentially capture more complex signals than modeled by hand-crafted features. For example, it is likely that features like curvature of streets, street signs, and traffic lights are also considered by MapLUR to estimate the air pollution. However, an analysis of these features remains future work. In particular, we believe that applying further model analysis techniques will enable researchers to find previously unknown features that can then be evaluated by experts and transferred to other, traditional land-use regression models. Thus, MapLUR's ability to extract interpretable features in combination with its inherent potential to model more complex relations of land-use and air pollution makes it a powerful tool for land-use regression.

Overfitting and generalization. Despite our strong results on the LAEI data set and gaining an intuition for how MapLUR works, we were not able to test generalizability across areas of interest due to the lack of similar data sets on different cities. While our setup allows for generalizability in principle, in practice, certain challenges may arise. In particular, deep-learning methods are prone to overfitting, i.e., they may underperform when applied to input data that is very dissimilar to or not covered by the training data [17]. However, we suspect that this problem is less pronounced for MapLUR, since map images are very structured and the model is therefore likely to see the vast majority of entity types that exist in the study area during training. Additionally, we want to stress again that previous models are often not applicable to new areas at all, since they often rely on data that is either not open or only locally available. Nonetheless, the model needs to be evaluated for every new application area before relying on its estimations. To use the model in new areas, it will usually be necessary to fit the model to some data from this area. Therefore, it is important to have training, validation, and test data sets that are similar in characteristics and representative of the whole study area. Typically, a simple random split is enough to achieve this [17].

Data Requirements and Application to Real World Data. MapLUR uses a CNN that contains a large number of weights due to its architecture. Training this deep-learning-based CNN requires more data than regular land-use regression models, which is why we have evaluated the approach on data from a model [2] instead of real world data. We have shown that MapLUR works well given the 3,000 data points of LAEI's modeled NO₂ concentrations. Since this number is far greater than most real world data sets for land-use regression, future work needs to investigate the possibility of applying models based on the DOG paradigm in more realistic settings. We believe that one promising approach for research in this direction is the use of transfer learning, which has been shown to be an effective way of dealing with low-resource settings in both the areas of computer vision [38] and natural language processing [28]. Transfer learning could be applied to MapLUR or other models based on the DOG paradigm by pre-training the model on a large data set, like for example the LAEI data, and then fine-tuning it to a smaller data set of real world measurements. The global nature of features used in models based on the DOG paradigm ensures that this approach is generally possible. Additional large-scale data sets can also be collected in the context of

mobile measuring campaigns [25, 41]. These data sets can then be used to provide further training data for the pre-training of DOG-based models.

Incorporating Distant Sources of Pollution. Our analysis has shown that using images depicting 80 m by 80 m areas for each data point leads to good results, as can be seen in Appendix B. However, previous approaches to air pollution estimation based on land-use regression have shown that it is useful to include information from wider surrounding areas in their features [49]. In the specific case of the MapLUR model, the size of the surrounding area that can be used is bounded by the resolution of the input images: If the area gets too large, then the resolution of 224×224 px is not sufficient to encode the corresponding image. While this could be countered by increasing the input resolution, this would significantly increase the cost of training and prediction. Therefore, it is an interesting direction for further research to develop models that can take into account larger surrounding areas without increasing the image resolution. This could for example be achieved by using stacked convolutional neural networks or a combination of convolutional and recurrent neural networks.

Integration of Additional Data Sources. While we have focused only on map images as input, MapLUR was able to outperform all considered baselines. Nevertheless, previous work has shown that additional information can greatly improve the performance of air pollution models. One example of such data would be elevation maps [49], which can be integrated into MapLUR in a way similar to Experiment 4, where we provided the CNN with additional map image layers. Beyond this, there is a wide variety of methods to provide deep-learning methods with additional information, which holds great potential to further improve our results [4, 39].

8 CONCLUSION

In this article, we have advocated DOG, a solely data-driven paradigm for air pollution estimation through land use regression. Models that follow this paradigm do not require manually engineered features and are based on data that is openly and globally available. This will ultimately result in models that are globally generalizable and can be applied in any area without modification. Working toward this goal, we have presented MapLUR, a deep-learning-based land-use regression model for air pollution estimation. We have shown that it can estimate NO_2 concentrations better than all considered baselines on a data set of modeled data from the Greater London area. While our analysis of MapLUR has shown that its data requirements are higher than commonly available data set sizes, we argued that transfer learning is a promising approach to alleviate this issue. We have also explored ways to analyze the factors that influence the prediction of this model, finding that a data-driven model architecture can be made interpretable by careful inspection of the trained model. Thus, overall, this article demonstrates the feasibility and advantages of our proposed data-driven paradigm DOG for land-use regression-based air pollution modeling.

Future directions encompass work to further reduce the data requirements of data-driven models, the development of a comprehensive framework for extracting and interpreting features, as well as in-depth studies on real-world data, large-scale mobile measurements, and different cities.

APPENDIXES

A SAMPLED LAEI CELLS

Figure 8 depicts the cells that we sampled from LAEI. The blue cells are used for training our models while the red cells are used to evaluate model performance for unseen locations.



Fig. 8. LAEI cells sampled for the experiments. The 3,000 blue cells are the training data set and the 1,500 red cells are the test data set. (Underlying OpenStreetMap Image: ©OpenStreetMap contributors).

B ANALYSIS OF AREA SIZE

In this work, we only used 80 m by 80 m images as inputs for our model. However, despite the potential evaluation issues with overlapping images described in Section 3.1, it is still interesting to see how our model behaves when it is able to see more or less of the surroundings. Therefore, the model was provided with OpenStreetMap images depicting square areas around the data point with side lengths of 60 m, 80 m, 100 m, 200 m, 500 m, and 1,000 m while maintaining a resolution of 224 px by 224 px. The mean results after 40 evaluations can be seen in Table 4.

Our model does not benefit from the increased image size as can be seen from both R^2 and RMSE. The mean performance decreases consistently with each increase in depicted area size over 80 m by 80 m. This implies that the potential evaluation issue with overlapping images, which is described in Section 3.1, is not very severe, since the model should be gaining performance with larger images otherwise. It also suggests that the very close surroundings are important and that information from further away is not helping. However, the model is suffering from performance loss with smaller areas than 80 m by 80 m. Thus, it seems that a side length of 80 m for the depicted square areas is optimal for MapLUR especially considering the fact that all other results are significantly different according to the Wilcoxon signed-rank test even after Bonferroni correction, since $p < \frac{0.05}{n}$ for $n = 9$.

Table 4. Evaluating the Influence of the Area Size Depicted in Each Map Image on MapLUR's Performance

Model	R ²	RMSE
Mean baseline	0.000	13.971
Linear regression	0.487	10.004
Multi-layer Perceptron	0.499	9.887
Random Forest	0.662	8.119
60 m	0.626	8.511
80 m	0.673	8.002
100 m	0.637	8.381
200 m	0.618	8.603
500 m	0.597	8.833
1000 m	0.390	10.856

The results suggest that 80 m by 80 m areas are optimal.

C EVALUATION METRICS

Given the desired target values $y = \{y_1, y_2, \dots, y_n\}$ and the model's output $\hat{y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n\}$, two commonly used metrics in land-use regression papers, namely, R² and root-mean-square error (RMSE) [14, 21, 37], are used to evaluate the model on the evaluation set.

On the one hand, R² describes how much of the target's variation is explained by the model:

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2},$$

where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the mean of all desired target values. The metric can take values from $-\infty$ to 1. A R² of 1 indicates a perfect fit. A value of 0 is achieved by always estimating the mean of the evaluation set's target values. Negative values indicate that the model is worse than always estimating the mean.

On the other hand, the RMSE is, as the name already suggests, the square root of the mean of the squared errors:

$$\text{RMSE}(y, \hat{y}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}.$$

Thus, RMSE can only take non-negative values, where 0 would be perfect for this metric and larger RMSEs are worse.

D MODEL BUILDING PROCEDURES FOR THE BASELINES

The following explains the model building procedures for the baseline methods in more detail.

D.1 Linear Regression

The first baseline model we consider is the commonly used linear regression. The model's development is based on a supervised stepwise selection that was used in the Escape project [20] for land-use regression model development before. Each predictor variable is ranked based on the model's adjusted R² from a univariate regression. The adjusted R² used by Eeftens et al. [20] is like the R² but penalizes adding variables that do not fit the model. Thus, ideally only independent variables that affect the dependent variable are used. If there are variables that are of the same

Table 5. Random Forest Hyperparameters Optimized Using Stochastic Search with the Corresponding Search Spaces

Hyperparameter	Search space
Number of trees	1 to 1,000
Fraction of features to consider at most per split	0.0 to 1.0
Minimum samples required to be a leaf node	1 to 100
Minimum samples required to split a node	2 to 20
Build trees with bootstrap samples	True or False

category but with different buffer sizes, then only the variable with the highest score is considered for use in the final model due to the high correlation of these variables between each other. The model starts with the variable that achieved the highest score. Thereafter, each one of the remaining variables is temporarily added to the model, evaluated, and the best performing variable is added to the model permanently if it increases the model's adjusted R^2 by at least 0.01. This is repeated until no variables are left. Then all selected variables with a p-value greater than 0.1 are removed and the resulting model is fit again, just like described in Eeftens et al. [20]. Finally, the variance inflation factors (VIFs) are calculated for each variable to quantify the increase in variance due to collinearity of the variables. If a variable has an VIF that is greater than 3, then the variable with the largest VIF is removed and the model is refit. In accordance with Eeftens et al. [20] this is also repeated until no variable has an VIF greater than 3.

D.2 Random Forest

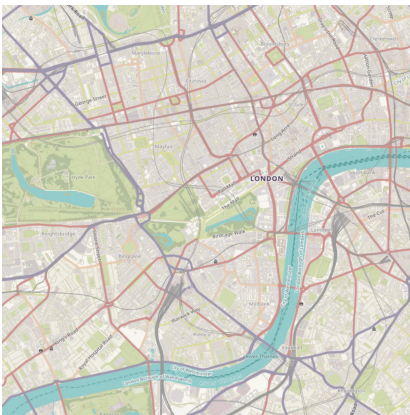
Another baseline is a Random Forest model that employs ensembles of decision trees for its estimations [9]. This model is built in a similar way to the procedure in Brokamp et al. [10]. The steps of the model building procedure are described in the following. First, the best buffer radii for each type of variable were determined based on the adjusted R^2 of a univariate regression on the training data set. As described in Section D.1 before, the adjusted R^2 is like the R^2 but penalizes adding variables that do not fit the model. Then, an initial Random Forest is considered for the following, which uses values that have shown to work decently in preliminary experiments: It builds 500 trees, considers half the available features when looking for the best split, and uses the default values of scikit-learn's Random Forest implementation for the other hyperparameters [48]. The best variables of each type are then fitted using this Random Forest to rank the variables based on the variable importance score of the Random Forest. Thereafter, the least important variables are removed iteratively. For each iteration, the model is fitted to the remaining variables on the training data set and the out of bag R^2 is calculated by estimating each training sample without using the trees that had the training sample in their bootstrap sample. This metric can be used with Random Forests to estimate performance without an independent test set. The set of variables that achieved the best performance is selected for further use. Then the best hyperparameters for the Random Forest are found by a stochastic search that samples hyperparameter values and evaluates them with a ten fold cross-validation on the folds of the training set. Each hyperparameter value is sampled from a uniform distribution with a specific search space. Table 5 shows these hyperparameters with their search spaces. This search ran for three hours and used six CPU cores of an Intel Xeon E5-2690V4 processor to fit as many models as possible. During this search, 1,051 sets of hyperparameters were evaluated. The best hyperparameters are then used to fit a Random Forest on the complete training set and evaluate the model on the test set.

D.3 Multi-layer Perceptron

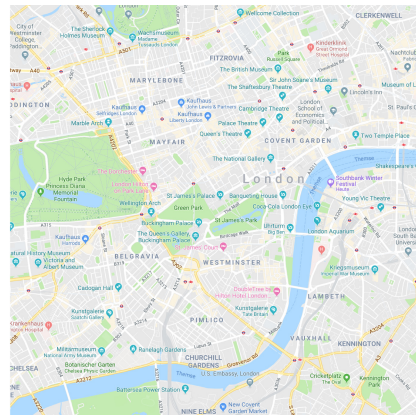
For the last baseline, we employ a multi-layer Perceptron, which is a type of neural network. For this model, we base our model building procedure on the one used by Alam and McNabola [3]. We adapt it slightly by first selecting the best buffer radius for each variable type based on the adjusted R^2 of an univariate regression. We use the adjusted R^2 again for this selection to be consistent with the model building procedures for linear regression and Random Forest described before. Then we search for the best performing architecture by evaluating models with different number of hidden layers and neurons for each layer with ten fold cross-validations on the training data set. This is done by randomly sampling the number of layers from a uniform distribution ranging from 1 to 3 hidden layers and randomly sampling the number of neurons for each layer from another uniform distribution that ranges from 1 to 30. These bounds are chosen, since they encompass the architectures of all previously published models that we found [1, 3, 12, 14, 37]. We randomly sample model architectures and evaluate them for 24 hours during which 2,629 architectures were evaluated. The best performing architecture is then trained on the complete training set and evaluated on the test set.

E COMPARISON OF OPENSTREETMAP AND GOOGLE MAPS IMAGES

Figure 9 compares Google Maps and OpenStreetMap images, showing the differences between both default styles shown on the services' websites. While OpenStreetMap provides at least four colors to denote different types of streets, Google Maps only uses two. Google Maps does not color-code all of the available information to ease the visual effort of the user. This, however, is not helpful for the CNN model, which works better with clear visual cues that denote entities.



(a) OpenStreetMap images (© OpenStreetMap contributors)



(b) Google Maps images (Map data: Google)

Fig. 9. Comparison of OpenStreetMap and Google Maps images. OpenStreetMap images contain more color-coded entities such as streets. As the visual style is highly important for the CNN to learn the task, the simpler Google Maps style produces worse results.

REFERENCES

- [1] Matthew Adams. 2015. *Advancing the use of mobile monitoring data for air pollution modelling*. Ph.D. Dissertation. McMaster University, Hamilton.

- [2] Air Quality Team (Greater London Authority). [n.d.]. London Atmospheric Emissions Inventory (LAEI). Retrieved from <https://data.london.gov.uk/dataset/london-atmospheric-emissions-inventory-2013>.
- [3] Md Saniul Alam and Aonghus McNabola. 2015. Exploring the modeling of spatiotemporal variations in ambient air pollution within the land use regression framework: Estimation of PM10 concentrations on a daily basis. *J. Air Waste Manage. Assoc.* 65, 5 (2015), 628–640.
- [4] Rahul Aralikkatte, Heather Lent, Ana Valeria Gonzalez, Daniel Hershovich, Chen Qiu, Anders Sandholm, Michael Ringaard, and Anders Søgaard. 2019. Rewarding Coreference Resolvers for Being Consistent with World Knowledge. Retrieved from <http://arxiv.org/abs/1909.02392>.
- [5] Yun Bai, Yong Li, Bo Zeng, Chuan Li, and Jin Zhang. 2019. Hourly PM2.5 concentration forecast using stacked autoencoder model with emphasis on seasonality. *J. Clean. Prod.* 224 (2019), 739–750.
- [6] Yun Bai, Bo Zeng, Chuan Li, and Jin Zhang. 2019. An ensemble long short-term memory neural network for hourly PM2.5 concentration forecasting. *Chemosphere* 222 (2019), 286–294.
- [7] Rob Beelen, Gerard Hoek, Danielle Vienneau, Marloes Eeftens, Konstantina Dimakopoulou, Xanthi Pedeli, Ming-Yi Tsai, Nino Künzli, Tamara Schikowski, Alessandro Marcon, et al. 2013. Development of NO₂ and NO_x land use regression models for estimating air pollution exposure in 36 study areas in Europe—The ESCAPE project. *Atmos. Environ.* 72 (2013), 10–23.
- [8] C. Bonferroni. 1936. Teoria statistica delle classi e calcolo delle probabilita. *Pubblazioni Del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze* 8 (1936), 3–62.
- [9] Leo Breiman. 2001. Random forests. *Mach. Learn.* 45, 1 (2001), 5–32.
- [10] Cole Brokamp, Roman Jandarov, M. B. Rao, Grace LeMasters, and Patrick Ryan. 2017. Exposure assessment models for elemental components of particulate matter in an urban environment: A comparison of regression and random forest approaches. *Atmos. Environ.* 151 (2017), 1–11.
- [11] Bert Brunekreef and Stephen T. Holgate. 2002. Air pollution and health. *Lancet* 360, 9341 (2002), 1233–1242.
- [12] Alexander G. Buevich, Alexander N. Medvedev, Alexander P. Sergeev, Dmitry A. Tarasov, Andrey V. Shichkin, Marina V. Sergeeva, and T. B. Atanasova. 2016. Modeling of surface dust concentrations using neural networks and kriging. In *AIP Conference Proceedings*, Vol. 1789. AIP Publishing, 020004.
- [13] David C. Carslaw and Sean D. Beevers. 2005. Estimations of road vehicle primary NO₂ exhaust emission fractions using monitoring data in London. *Atmos. Environ.* 39, 1 (2005), 167–177.
- [14] Alexandre Champendal, Mikhail Kanevski, and Pierre-Emmanuel Huguenot. 2014. Air pollution mapping using non-linear land use regression models. In *Proceedings of the International Conference on Computational Science and Its Applications*. Springer, 682–690.
- [15] Ralph B. d’Agostino. 1971. An omnibus test of normality for moderate and large size samples. *Biometrika* 58, 2 (1971), 341–348.
- [16] Ralph B. d’Agostino and Egon S. Pearson. 1973. Tests for departure from normality. Empirical results for the distributions of b_2 and $\sqrt{b_1}$. *Biometrika* 60, 3 (1973), 613–622.
- [17] Howard B. Demuth, Mark H. Beale, Orlando De Jess, and Martin T. Hagan. 2014. *Neural Network Design* (2nd ed.). Martin Hagan, Stillwater, OK.
- [18] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR’09)*.
- [19] Harris Drucker, Christopher J. C. Burges, Linda Kaufman, Alex J. Smola, and Vladimir Vapnik. 1997. Support vector regression machines. In *Advances in Neural Information Processing Systems*. MIT Press, 155–161.
- [20] Marloes Eeftens, Rob Beelen, Kees de Hoogh, Tom Bellander, Giulia Cesaroni, Marta Cirach, Christophe Declercq, Audrius Dedele, Evi Dons, Audrey de Nazelle et al. 2012. Development of land use regression models for PM2.5, PM2.5 absorbance, PM10 and PMcoarse in 20 European study areas; results of the ESCAPE project. *Environ. Sci. Technol.* 46, 20 (2012), 11195–11205.
- [21] Marloes Eeftens, Ming-Yi Tsai, Christophe Ampe, Bernhard Anwander, Rob Beelen, Tom Bellander, Giulia Cesaroni, Marta Cirach, Josef Cyrus, Kees de Hoogh, et al. 2012. Spatial variation of PM2.5, PM10, PM2.5 absorbance and PMcoarse concentrations between and within 20 European study areas and the relationship with NO₂—Results of the ESCAPE project. *Atmos. Environ.* 62 (2012), 303–317.
- [22] Derek M. Elsom. 1992. *Atmospheric Pollution: A Global Problem*. Blackwell, Oxford.
- [23] Junxiang Fan, Qi Li, Junxiong Hou, Xiao Feng, Hamed Karimian, and Shaofu Lin. 2017. A spatiotemporal prediction framework for air pollution based on deep RNN. *ISPRS Ann. Photogram. Remote Sens. Spatial Info. Sci.* 4 (2017), 15.
- [24] Google LLC. 2018. Google Maps. <https://maps.google.com>.
- [25] David Hasenfrazt, Olga Saukh, Christoph Walsler, Christoph Hueglin, Martin Fierz, and Lothar Thiele. 2014. Pushing the spatio-temporal resolution limit of urban air pollution maps. In *Proceedings of the IEEE International Conference on Pervasive Computing and Communications (PerCom’14)*. IEEE, 69–77.

- [26] Gerard Hoek, Rob Beelen, Kees de Hoogh, Danielle Vienneau, John Gulliver, Paul Fischer, and David Briggs. 2008. A review of land-use regression models to assess spatial variation of outdoor air pollution. *Atmos. Environ.* 42, 33 (2008), 7561–7578.
- [27] Gerard Hoek, Marloes Eeftens, Rob Beelen, Paul Fischer, Bert Brunekreef, K. Folkert Boersma, and Pepijn Veeffkind. 2015. Satellite NO₂ data improve national land use regression models for ambient NO₂ in a small densely populated country. *Atmos. Environ.* 105 (2015), 173–180. DOI: <https://doi.org/10.1016/j.atmosenv.2015.01.053>
- [28] Jeremy Howard and Sebastian Ruder. 2018. Universal Language Model Fine-tuning for Text Classification. Retrieved from <http://arxiv.org/abs/1801.06146>.
- [29] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- [30] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [31] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. 2017. Self-normalizing neural networks. In *Advances in Neural Information Processing Systems*. MIT Press, 971–980.
- [32] Stéphane Lathuilière, Pablo Mesejo, Xavier Alameda-Pineda, and Radu Horaud. 2019. A comprehensive analysis of deep regression. *IEEE Trans. Pattern Anal. Mach. Intell.* (2019).
- [33] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* 521, 7553 (2015), 436.
- [34] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [35] Xiaoli Li, Aorong Luo, Jiangeng Li, and Yang Li. 2019. Air pollutant concentration forecast based on support vector regression and quantum-behaved particle swarm optimization. *Environ. Model. Assess.* 24, 2 (Apr. 2019), 205–222. DOI: <https://doi.org/10.1007/s10666-018-9633-3>
- [36] Yuncheng Li, Jifei Huang, and Jiebo Luo. 2015. Using user generated online photos to estimate and monitor air pollution in major cities. In *Proceedings of the 7th International Conference on Internet Multimedia Computing and Service*. ACM, 79.
- [37] Wu Liu, Xiaodong Li, Zuo Chen, Guangming Zeng, Tomás León, Jie Liang, Guohe Huang, Zhihua Gao, Sheng Jiao, Xiaoxiao He, et al. 2015. Land use regression models coupled with meteorology to model spatial and temporal variability of NO₂ and PM10 in Changsha, China. *Atmos. Environ.* 116 (2015), 272–280.
- [38] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. 2018. Exploring the Limits of Weakly Supervised Pretraining. Retrieved from <http://arxiv.org/abs/1805.00932>.
- [39] Julian McAuley and Jure Leskovec. 2012. Image labeling on a network: Using social-network metadata for image classification. In *Proceedings of the European Conference on Computer Vision*. Springer, 828–841.
- [40] Xia Meng, Li Chen, Jing Cai, Bin Zou, Chang-Fu Wu, Qingyan Fu, Yan Zhang, Yang Liu, and Haidong Kan. 2015. A land use regression model for estimating the NO₂ concentration in Shanghai, China. *Environ. Res.* 137 (2015), 308–315.
- [41] Denise R. Montagne, Gerard Hoek, Jochem O. Klompmaker, Meng Wang, Kees Meliefste, and Bert Brunekreef. 2015. Land use regression models for ultrafine particles and black carbon based on short-term monitoring predict past spatial variation. *Environ. Sci. Technol.* 49, 14 (2015), 8712–8720.
- [42] David W. Morley and John Gulliver. 2018. A land use regression variable generation, modelling and prediction tool for air pollution exposure assessment. *Environ. Model. Softw.* 105 (2018), 17–23.
- [43] Sheena Mutttoo, Lisa Ramsay, Bert Brunekreef, Rob Beelen, Kees Meliefste, and Rajen N. Naidoo. 2018. Land use regression modelling estimating nitrogen oxides exposure in industrial south Durban, South Africa. *Sci. Total Environ.* 610 (2018), 1439–1447.
- [44] Vinod Nair and Geoffrey E. Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML'10)*. 807–814.
- [45] OpenStreetMap contributors. 2018. Apache Module Mod_tile. Retrieved from https://github.com/openstreetmap/mod_tile.
- [46] OpenStreetMap contributors. 2018. Planet Dump Retrieved from Retrieved from <https://planet.osm.org>; <https://www.openstreetmap.org>.
- [47] E. G. Ortiz-García, S. Salcedo-Sanz, Á. M. Pérez-Bellido, J. A. Portilla-Figueras, and L. Prieto. 2010. Prediction of hourly O₃ concentrations using support vector regression algorithms. *Atmos. Environ.* 44, 35 (2010), 4481–4488.
- [48] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* 12 (2011), 2825–2830.
- [49] Patrick H. Ryan and Grace K. LeMasters. 2007. A review of land-use regression models for characterizing intraurban air pollution exposure. *Inhal. Toxicol.* 19 (2007), 127–133.
- [50] Vikas Singh. 2016. Higher pollution episode detection using image classification techniques. *Environ. Model. & Assess.* 21, 5 (Oct. 2016), 591–601. DOI: <https://doi.org/10.1007/s10666-015-9497-8>

- [51] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. 2014. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*.
- [52] Meng Wang, Rob Beelen, Tom Bellander, Matthias Birk, Giulia Cesaroni, Marta Cirach, Josef Cyrys, Kees de Hoogh, Christophe Declercq, Konstantina Dimakopoulou, et al. 2014. Performance of multi-city land use regression models for nitrogen dioxide and fine particles. *Environ. Health Perspect.* 122, 8 (2014), 843.
- [53] Kathrin Wolf, Josef Cyrys, Tatiana Hrciniková, Jianwei Gu, Thomas Kusch, Regina Hampel, Alexandra Schneider, and Annette Peters. 2017. Land use regression modeling of ultrafine particles, ozone, nitrogen oxides and markers of particulate matter pollution in Augsburg, Germany. *Sci. Total Environ.* 579 (2017), 1531–1540.
- [54] Jiansheng Wu, Jiacheng Li, Jian Peng, Weifeng Li, Guang Xu, and Chengcheng Dong. 2015. Applying land use regression model to estimate spatial variation of PM_{2.5} in Beijing, China. *Environ. Sci. Pollution Res.* 22, 9 (2015), 7045–7061.
- [55] Jingjing Xie, Xiaoxue Wang, Yu Liu, and Yun Bai. 2018. Autoencoder-based deep belief regression network for air particulate matter concentration forecasting. *J. Intell. Fuzzy Syst.* 34, 6 (2018), 3475–3486.
- [56] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*.
- [57] Chao Zhang, Junchi Yan, Changsheng Li, Hao Wu, and Rongfang Bie. 2018. End-to-end learning for image-based air quality level estimation. *Mach. Vision Appl.* 29, 4 (2018), 601–615.

Received May 2019; revised September 2019; accepted January 2020