

FaceLift: A transparent deep learning framework recreating the urban spaces people intuitively love

SAGAR JOGLEKAR, King's College London, Department of Informatics

DANIELE QUERCIA, Nokia Bell Labs, Cambridge, UK

MIRIAM REDI, Nokia Bell Labs, Cambridge, UK

LUCA MARIA AIELLO, Nokia Bell Labs, Cambridge, UK

TOBIAS KAUER, Nokia Bell Labs, Cambridge, UK

NISHANTH SASTRY, King's College London, Department of Informatics

There has been an explosive growth of deep learning technologies and their competency in the recent years, resulting in cross disciplinary use cases of deep learning enabled tools. In the area of computer vision and urban informatics, deep learning techniques have recently been used to predict whether urban scenes are likely to be considered beautiful, and it turns out that these techniques do so quite accurately. However, the technology falls short when it comes to generating actionable insights for AI assisted urban design. To support urban interventions, one needs to go beyond *predicting* beauty, and tackle the challenge of *recreating* beauty and *explaining* the predictors of beauty. And for these explanations to be of any use to the target audience of such tools, they need to be grounded in the literature and language of the target users. Unfortunately, deep learning techniques have not been designed with that challenge in mind. Given their “black-box nature”, these models cannot be directly used to explain why a particular urban scene is deemed to be beautiful. To partly fix that, we propose a deep learning framework (which we name FaceLift) that is able to both *beautify* existing Google Street views and *explain* which urban elements make those transformed scenes beautiful, in the vocabulary of urban design science. To quantitatively evaluate our framework, we cannot resort to any existing metric (as the research problem at hand has never been faced before) and need to formulate new ones. These new metrics should ideally capture the presence (or absence) of elements that make urban spaces great. They ideally should also be computable using current computer vision techniques. Upon a review of the urban planning literature, we identify five main metrics: walkability, green spaces, openness, landmarks and visual complexity. For all the five metrics, the beautified scenes meet the expectations set by the literature on what great spaces tend to be made of. The transformations and their explanations are also found to be very helpful in understanding interventions for beautification, which we validate using a 20-participant expert survey. These results suggest that, in the future, as our framework’s components are further researched and become better and more sophisticated, it is not hard to imagine technologies that will be able to accurately and efficiently support architects and planners in the design of the spaces we intuitively love.

Additional Key Words and Phrases: Deep learning, Generative networks, Urban Beauty, Computer Vision

ACM Reference format:

Sagar Joglekar, Daniele Quercia, Miriam Redi, Luca Maria Aiello, Tobias Kauer, and Nishanth Sastry. 2017. FaceLift: A transparent deep learning framework recreating the urban spaces people intuitively love. 1, 1, Article 1 (January 2017), 19 pages.

DOI: 10.1145/nnnnnnn.nnnnnnn

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 ACM 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.

© 2017 ACM. XXXX-XXXX/2017/1-ART1 \$15.00

DOI: 10.1145/nnnnnnn.nnnnnnn

1 INTRODUCTION

Whether a street is considered beautiful is subjective, yet research has shown that there are specific urban elements that are universally considered beautiful: from greenery, to small streets, to memorable spaces [1, 32, 35]. These elements are those that contribute to the creation of what the urban sociologist Jane Jacobs called ‘urban vitality’ [16].

Given that, it comes as no surprise that computer vision techniques can automatically analyse pictures of urban scenes and accurately determine the extent to which these scenes are, *on average*, considered beautiful. Deep learning has greatly contributed to increase these techniques’ accuracy [10].

However, urban planners and architects are interested in urban interventions and, as such, they wish to go beyond technologies that are only able to predict beauty scores. These interests stem from the fact that the spaces we live in can be linked with several aspects of human life such as mental health[37], inequality [35] or cultural shifts [15]. They often called for technologies that would make easier to recreate beauty in urban design [7]. Deep learning, by itself, is not fit for purpose. It is not meant to recreate beautiful scenes, not least because it cannot provide any explanation on why a scene is deemed beautiful, or which urban elements are predictors of its beauty.

To partly fix that, we propose a deep learning framework (which we name FaceLift) that is able to both *generate* a beautiful scene (or, better, *beautify* an existing scene) and *explain* why that scene is beautiful. This opens up a possibility of using technology in urban planning efforts like decision making based of subjective opinions, participatory urban planning and promotion of restorative urban design such as green spaces and walkable areas. Through this work, we make two main contributions:

- We propose a deep learning framework that is able to learn whether a particular set of Google Street views are beautiful or not, and based on that training, is able to both *beautify* existing views which were deemed not to be as beautiful and *explain* which urban elements make these views beautiful (Section 3). The explanations come in the form of predictors of urban beauty, measured using computer vision tools.
- We quantitatively evaluate whether the framework is able to actually produce beautified scenes (Section 4). We do so by proposing a family of five urban design metrics that we have formulated based on a thorough review of the literature in urban planning. For all these five metrics, the framework passes with flying colors: with minimal interventions, beautified scenes are twice as walkable as the original scenes, for example. Also, after building an interactive tool with “FaceLifted” scenes in Boston and presenting it to twenty experts in architecture, we found that the majority of them agreed on three main areas of our work’s impact: decision making, participatory urbanism, and promotion of restorative spaces among the general public.

For sake of brevity, we will use the term ‘Urban Scene’ through out the paper to address an arbitrary Google Street View image. The image is fetched from a particular latitude and longitude point on the map. In the rest of the paper we explore related literature across various tracks of urban perceptions and urban beauty in Section 2. We then describe in detail the Facelift framework in Section 3. The evaluation of the framework is described in detail in Section 4. We conclude by pointing out some limitations and biases that might well guide future work (Section 5).

2 RELATED WORK

Previous work has focused on collecting ground truth data about how people perceive urban spaces, on predicting urban qualities from visual data, and on generating synthetic images that enhance a given quality (e.g., beauty).

Perception of physical spaces. The literature in the area of quantifying people perception of urban environments is pretty rich. From the seminal work about urban vitality by Jane Jacobs [16] to imagining urban design through patterns [1], there has been a continuous effort to understand and intervene to make our cities more liveable and enjoyable. A lot of work was done in the field of human behaviour analysis, e.g., Roger Ulrich's work to understand affective responses to urban environments [43]. Some studies looked into self rated perception of urban aesthetics against well known metrics like complexity [18] or perception of nature [17]. Other empirical works done used trained survey takers to understand how people perceive scenic beauty [34]. All this work in the fields of psychology, environmental design and behavioural sciences showed that humans' perception of physical spaces is quite predictable, which makes use of technologies like deep learning more relevant in this area.

Ground truth of urban perceptions. So far, the most detailed studies of perceptions of urban environments and their visual appearance have relied on personal interviews and observation of city streets: for example, some researchers relied on annotations of video recordings by experts [36], while others have used participant ratings of simulated (rather than existing) street scenes [23]. The Web has recently been used to survey a large number of individuals. Place Pulse is a website that asks a series of binary perception questions (such as 'Which place looks safer [between the two?]') across a large number of geo-tagged images [35]. In a similar way, Quercia *et al.* collected pairwise judgments about the extent to which urban scenes are considered quiet, beautiful and happy [32] to then recommend pleasant paths in the city [33]. They were then able to analyze the scenes together with their ratings using image-processing tools, and found that the amount of greenery in any given scene was associated with all three attributes and that cars and fortress-like buildings were associated with sadness. Taken all together, their results pointed in the same direction: urban elements that hinder social interactions were undesirable, while elements that increase interactions were the ones that should be integrated by urban planners to retrofit cities for greater happiness. Some studies also linked aesthetics of physical spaces to physical activity: Ball *et al.* [6] collected 3.3k self reported surveys and showed that urban aesthetics have a positive effect on the urge of walking. In another work [12], Giles *et al.* interviewed 1.8k participants and found that attractiveness of public open spaces impact positively on increased walking. Finally, a large scale study using a crowd sourced interface¹ looked at relationship between 'scenic-ness' of a place with land cover and geography.

Deep learning and the city. Computer vision techniques have increasingly become more sophisticated. Deep learning techniques, in particular, have been recently used to accurately predict urban beauty [10, 38], urban change [28], and even crime [3, 8]. These works also did some interesting analysis of the data to understand how safety, depression, beauty and other such dimensions are perceived across urban spaces. [10] also utilized deep learning methods to train models capable of comparing two urban images for their perception values in terms of beauty et.al. Recent work has also explored utility of deep learning techniques in understanding relationship between urban frontage and housing prices [21]. Another work looked at predicting house prices using deep learning on satellite as well as street view images [20]. With the advent of augmented reality, the application of GANs to generate urban objects, so as to simulate urban driving scenes have also been explored [2]. This shows that GANs can be immensely useful in problems where you need to model real world and generate samples that mimic the real world as close as possible. While these works, similar to ours, model aspects of urban perception, they do not dive into the reasoning aspect behind these models.

Generative models. Since the introduction of Generative adversarial Networks [13], deep learning has recently been used not only to analyse existing images but also to generate new ones that mimic

¹<http://scenic.mysociety.org/>

Symbol	Meaning
I_i	Original urban scene
Y	Set of annotation classes for urban scenes (e.g., beautiful, ugly)
y_i	Annotation class in Y (e.g., beautiful)
\hat{I}_j	Template scene (synthetic image)
I'	Target Image
C	Beauty Classifier

Table 1. Notations

certain properties of training data. This family of deep networks has evolved into various forms, from super resolution image generators [22], to networks that could in-paint from context [30]. In the past couple of years, there have been papers which exploit generative version of neural nets to delve into the aspects of explainability. Recently GANs were used to do semantic segmentation of images [24]. This approach paves a way for using latent knowledge learned by the classifiers to explain semantics in the image. Similar approaches have been used to generate images conditioned on specific visual attributes [45] or generation of faces [42] or images of whole people [26]. Nguyen *et al.* [29] used generative networks to create a natural-looking image that maximizes a specific neuron. This method was used to bring out the latent representation of an image, that maximizes its probability of a particular class. In theory, the resulting image is the one that “best activates” the neuron under consideration. In practice, it is still a synthetic template that needs further processing to look realistic.

To sum up, a lot of work has gone into collecting ground truth data about how people tend to perceive urban spaces, and into building accurate predictions models of urban qualities. However, little work has gone into models that generate realistic urban scenes and that offer human-interpretable explanations of what they generate.

3 FACELIFT FRAMEWORK

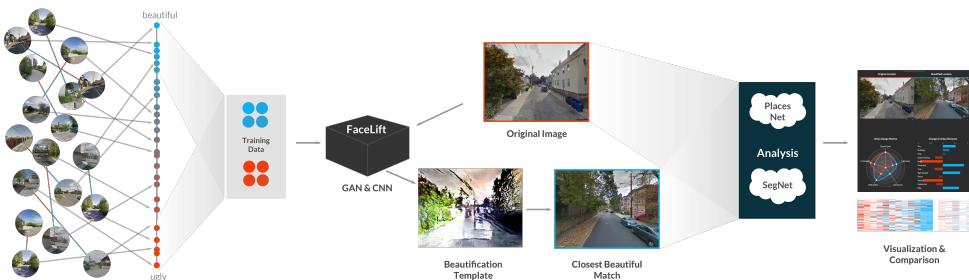


Fig. 1. An illustration of the FaceLift framework.

The goal of FaceLift is to take as input a geo-located urban scene and give as output its transformed (beautified) version. To that end, it performs in five steps:

- **Curating urban scenes** It is common knowledge that deep learning systems need immense amount of data. In this first step we try to develop a sound framework for curating and augmenting annotated images, on which the model could be trained.
- **Training a beauty classifier** To generate beauty, you first need a reliable model that could learn the representation of beauty. To achieve this, we train a deep learning model that could distinguish beautiful urban scenes from non-beautiful urban scenes.
- **Generating a synthetic beautified scene** Based on the learned representation of beauty, we train a generative model which could augment the beauty of an input urban scene.
- **Retrieving a realistic beautified scene** as showcased in Figure 1, the generated images are representations of beautified input urban scene in a latent space. This latent representation needs to be transformed back to a realistic looking image, using retrieval.
- **Identifying the urban elements characterizing the beautified scene** In the final step, the framework explains changes introduced in the transformation process in terms of literature-driven urban design metrics, and quantifies these changes as metrics for urban beauty.

Step 1 Curating Urban Scenes

To begin with, we need highly curated training data with labels reflecting urban beauty. We start with the Place Pulse dataset that contains 100k Google Street Views across 56 cities around the world [10]. These scenes are labeled in terms of whether the corresponding places are likely to be perceived beautiful, depressing, rich, and safe. We focus only on those scenes that are labeled in terms of beauty and that have at least three judgments. This leave us with roughly 20,000 scenes. To transform judgments into beauty scores, we use the TrueSkill algorithm [14], which gives us a way of partitioning the scenes into two sets (Figure 2): one containing beautiful scenes, and the other containing ugly scenes. The resulting set of scenes is too small for training any deep learning model without avoiding over-fitting though. As such, we need to augment such a set.

We do so in two ways. First, we feed each scene’s location into the Google Streetview API to obtain the snapshots of the same location at different camera angles (i.e., at $\theta \in -30^\circ, -15^\circ, 15^\circ, 30^\circ$). However, the resulting dataset is still too small for robust training. Therefore, again, we feed each scene’s location into the Google Streetview API, but now we do so to obtain other scenes at distance $d \in \{10, 20, 40, 60\}$ meters. This will greatly expand our set of scenes, but it might do so at the price of introducing scenes whose beauty scores have little to do with the original scene’s. To fix that, we take only the scenes that are *similar* to the original one (we call this way of augmenting “conservative translation”). To compute the similarity between a pair of scenes, we represent the two scenes with visual features derived from the FC7 layer of PlacesNet and compute the similarity between the two corresponding feature vectors [46]. For all scenes at increasing distance $d \in \{10, 20, 40, 60\}$ meters, we take only those whose similarity scores with the original scene is above a threshold. In a conservative fashion, we choose that threshold to be the median similarity between rotated and original scenes (those of the first augmentation step).

To make sure this additional augmentation has not introduced any unwanted noise, we consider two sets of scenes: one containing those that have been taken during this last step, i.e. the one with high similarity to the original scenes (*taken-set*), and the other containing those that have been filtered away (*filtered-set*). Each scene is then scored with PlacesNet [46] and is represented with the five most confident scene labels. We then aggregate labels at set level, by computing each label’s frequency on the *taken-set* and on the *filtered-set*. Finally, we characterize each label’s propensity to be correctly augmented as: $\text{prone}(\text{label}) = \text{fr}(\text{label}, \text{taken-set}) - \text{fr}(\text{label}, \text{filtered-set})$. This reflects the extent to which a scene with a given label is prone to be augmented or not. From Figure 4, we find that, as one would expect, scenes that contain highways, fields and bridges can be augmented at

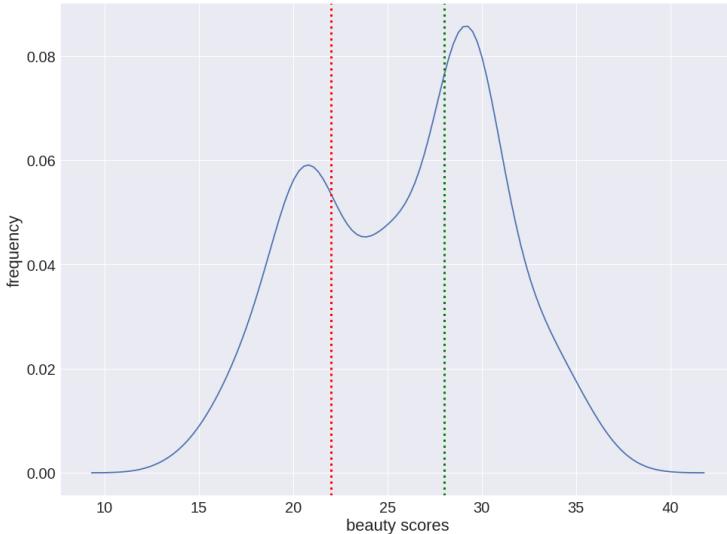


Fig. 2. Frequency distribution of beauty scores. The red and green lines represent the thresholds below and above which images are considered ugly and beautiful. Conservatively, images in between are discarded.

Augmentation	Accuracy (Percentage)
None	63
Rotation	68
Rotation + Translation	64
Rotation + Conservative Translation	73.5

Table 2. Percentage accuracy for our beauty classifier trained on differently augmented sets of urban scenes.

increasing distances while still showing resemblances to the original scene; by contrast, scenes that contain gardens, residential neighborhoods, plazas, and skyscrapers cannot be easily augmented, as they are often found in high density parts of the city, where there is tremendous diversity within short distances.

Step 2 Training a beauty classifier

Having this highly curated set of labeled urban scenes, we are now ready to train classifier C with labels reflecting our beauty assessments. We use the CaffeNet architecture, a modified version of AlexNet [19, 41]. The classifier network has a conventional convolutional network architecture, with 5 convolutional layers, interleaved with Max pooling and normalization layers. The network is terminated with two 4096 dimensional fully connected layers interleaved with dropout [40] layers, with dropout ratio of 0.5, to prevent over-fitting of the network. Finally the classifier categorizes the images into one of two classes (beautiful(1), ugly(0)) through a Softmax layer that computes probabilities of class membership.

The training is done on a 70% split of the data, and the testing on the remaining 30%. All this is done on increasingly augmented sets of data. We start from our 20k images and progressively augment them with the snapshots obtained with the 5-angle camera rotations, and then with the exploration of scenes at increasing distance $d \in \{10, 20, 40, 60\}$ meters. The idea behind data augmentation is that

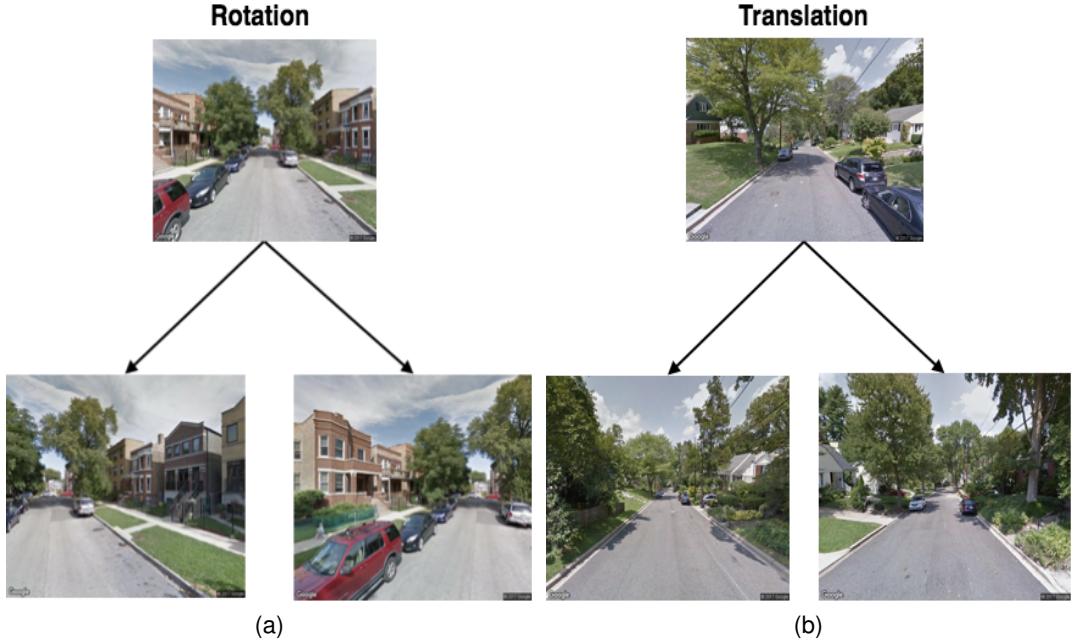


Fig. 3. Two types of augmentation: (a) rotation of the Street Views camera (based on rotation); and (b) exploration of scenes at increasing distances (based on translation).

accuracy would increase with it. Indeed it does (Table 2): it goes from 63% on the set of original scenes to as much as 73.5% on the set of fully augmented scenes, which is a notable increase in accuracy for such classes of classification tasks. As a baseline, we compare with the models trained by Dubey et.al [10] on the same seed data that we use for our pipeline. They report that their models perform at 70% accuracy in the task of picking a beautiful image amongst any two given images. Albeit the set-up of our model is not to compare two images but just to classify a particular image in a binary class, this baseline shows that our model is showing a comparable performance in beauty classification.

Step 3 Generating a synthetic beautified scene

Having this trained classifier at hand, we can then build a generator of synthetic beautified scenes. There are two components of this step. The first is a generator of synthetic scenes given an a-priori feature vector f , extracted from image I_f . For this component, we retrain the GAN described by Dosovitskiy and Brox [9] on the curated urban scene dataset described in Section 3. This network is trained by optimizing along the principles of Generative Adversarial Networks [13] i.e., maximizing the confusion for the discriminator between generated image $G(f)$ and real dataset images I_f . The resulting generator is able to reproduce synthetic urban scenes to an impressive level of detail. Figure 3 shows a few examples of how the generator performs by comparing a real world image I_f , and the generated image $G(f)$. The second component is a concatenation of the trained generator with the beauty classifier described in Section 3, as shown in Figure 5. This results in the end to end model that, given the two classes ugly y_i and beautiful y_j , transforms any original scene I_i of class y_i (e.g., ugly scene) into template scene \hat{I}_j that maximizes class y_j (e.g., beautified template scene).

More specifically, given an input image I_i known to be of class y_i (e.g., ugly), our technique outputs \hat{I}_j , which is a more beautiful version of it (e.g., I_i is morphed towards the average representation of

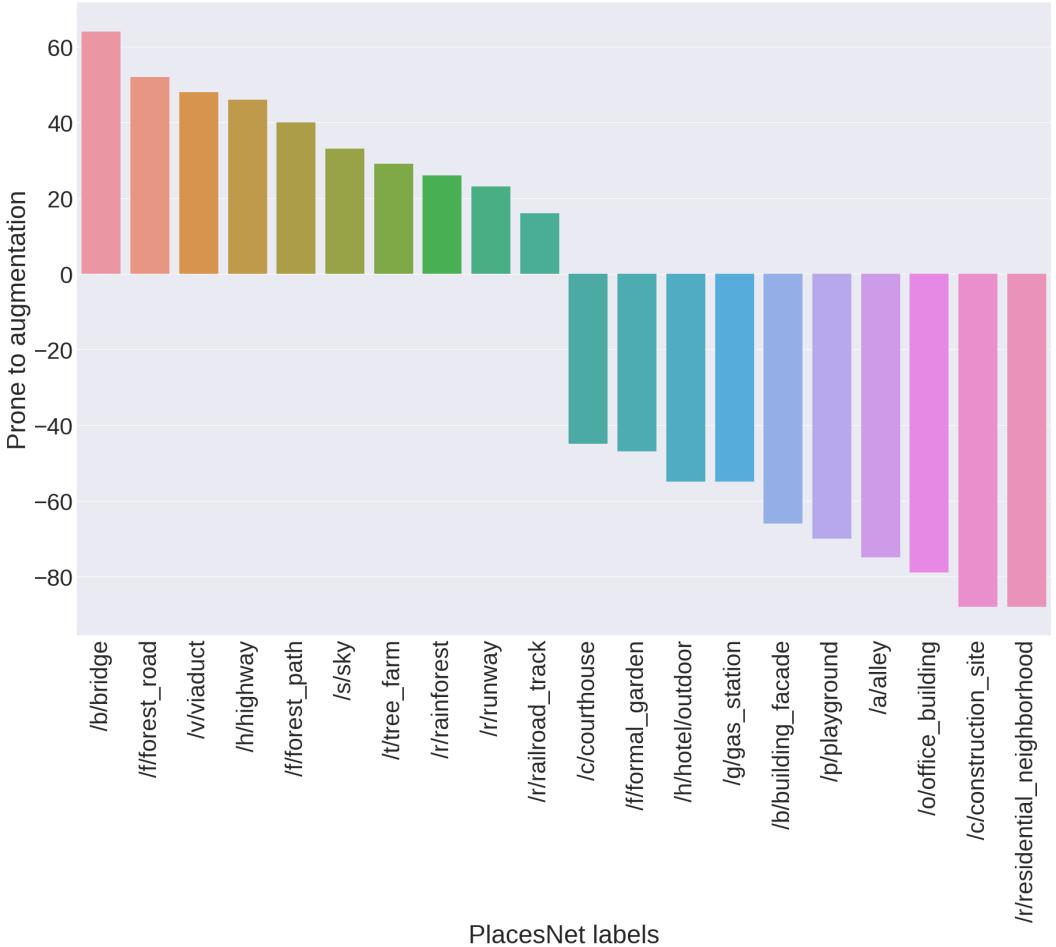


Fig. 4. The types of scene that have greater propensity to be correctly augmented with similar scenes at increasing distances.

a beautiful scene) while preserving I_i 's details. The technique does so using the “Deep Generator Network for Activation Maximization” (*DGN-AM*) [29]. Given an input image I_i , *DGN-AM* iteratively re-calculates the color of I_i 's pixels in a way the output image \hat{I}_j both maximizes the activation of neuron y_j (e.g., the “beauty neuron”) and looks “photo realistic”, which is done by conditioning the maximization to an “image prior”. This is equivalent to finding the feature vector f that maximizes the following expression:

$$\hat{I}_j = G(f) : \arg \max_f (C_j(G(f)) - \lambda \|f\|) \quad (1)$$

where:

- $G(f)$ is the image synthetically generated from the candidate feature vector f ;
- $C_j(G(f))$ is the activation value of neuron y_j in the scene classifier C (the value to be maximized);

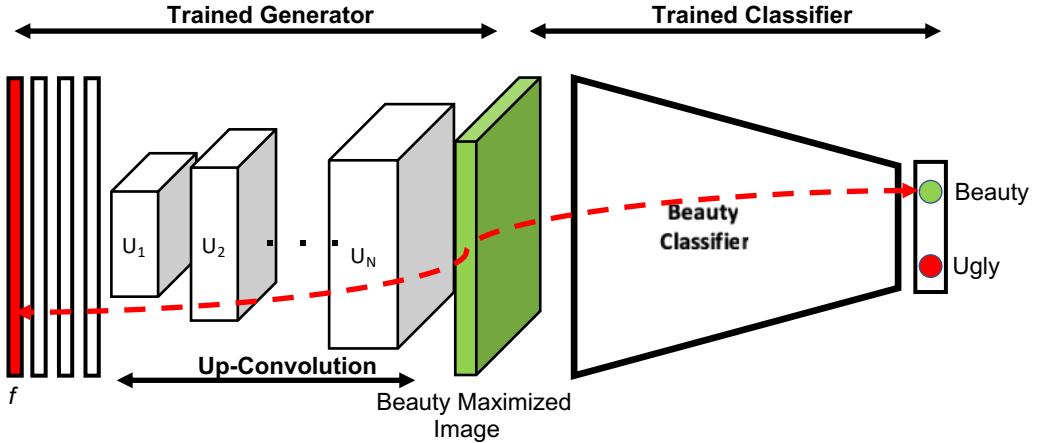


Fig. 5. The figure shows the architecture of the synthetic beauty generation pipeline. The pipeline is a Generator trained on all the curated Streetview images in the dataset, cascaded with the trained classifier described in Section 3. The green block represents the beauty maximized template \hat{I}_j . The red arrow describes the forward and backward pass that is optimized on.

- λ is a L_2 regularization term.

Here the initialization of f is key. If f were to be initialized with random noise, then the resulting $G(f)$ would be the average representation of category y_j (of, e.g., beauty). Instead, since f is initialized with the feature vector corresponding to I_i , then the resulting maximized $G(f)$ is I_i 's version “morphed to become more beautiful”. Maximizing beauty of an already beautiful urban scene, would yield in a saturated template \hat{I}_j . For this reason, to generate an image that maximizes the beauty neuron in the classifier C , we supply the network with an a-priori image that most definitely lies in the class y_i . The constraint is reverse for maximizing for class y_i .

Step 4 Returning a realistic beautified scene

We now have template scene \hat{I}_j (which is a synthetic beautified version of original scene I_i) and need to retrieve a realistic looking version of it. We do so by: *i*) representing each of our original scenes in Step 1 (including \hat{I}_j) as a 4096 dimensional feature vector derived from the FC7 layer of the PlacesNet [46]; *ii*) computing the distance (as L_2 Norm) between \hat{I}_j 's feature vector and each of the original scene's feature vector; and *iii*) selecting the original scene most similar (smaller distance) to \hat{I}_j . This results into the selection of the beautified scene I_j .

Step 5 Identifying characterizing urban elements

Since original scene I_i and beautified scene I_j are real scenes and we make sure that they maintain the same structural characteristics (e.g., point of view, layout), we can easily compare them in terms of presence or absence of SegNet's and PlacesNet's labels. That is, we can determine how the original scene and its beautified version differ in terms of urban design elements. This step required us to develop metrics inspired from urban design literature, to quantify the changes in elements. A detailed description of the characterization and evaluation would follow in Section 4.

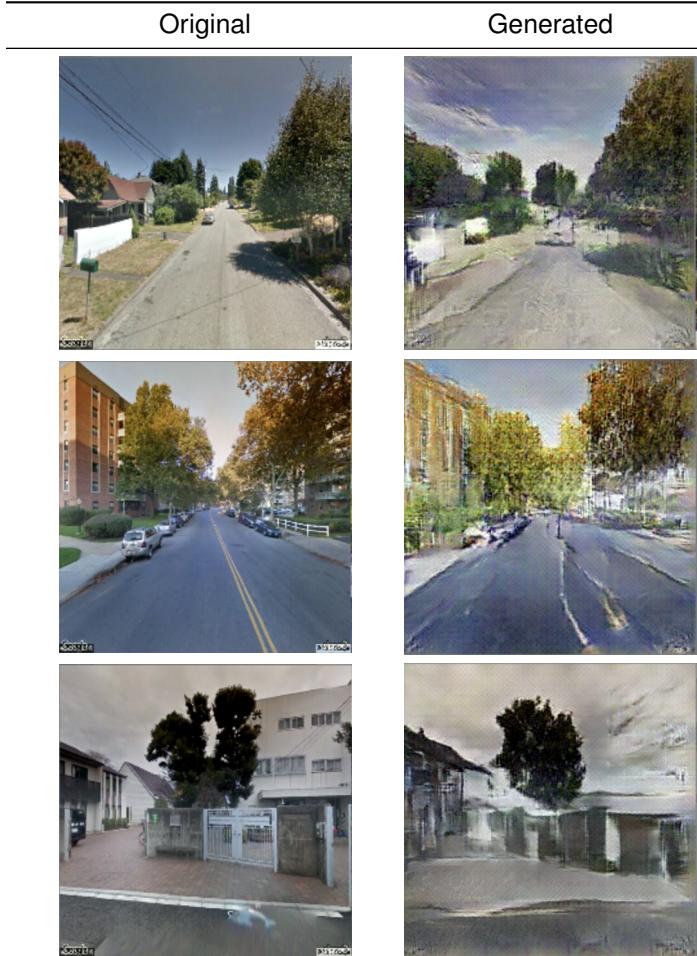


Table 3. Sample examples of the outputs of the Generator network. The original scenes and the generated scenes are shown side by side.

4 EVALUATION

The goal of FaceLift is to transform existing urban scenes into versions that: *i*) people perceive more beautiful; *ii*) contain urban elements typical of great urban spaces; *iii*) are easy to interpret; and *iv*) architects and urban planners find useful. To ascertain whether FaceLift meets that composite goal, we answer the following questions next:

- Q1** Do individuals perceive “FaceLifted” scenes to be beautiful?
- Q2** Does our framework produce scenes that possess urban elements typical of great spaces?
- Q3** Which urban elements are mostly associated with beautiful scenes?
- Q4** Do architects and urban planners find FaceLift useful?

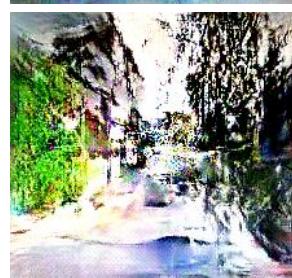
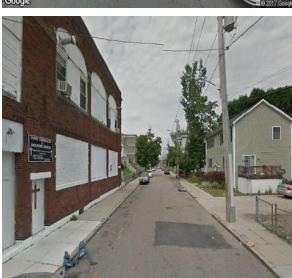
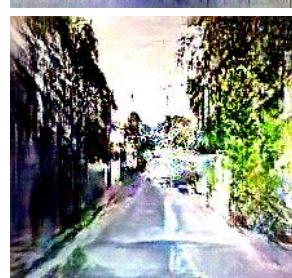
Original (I_i)	Latent Beauty representation (\hat{I}_j)	Beautified (I_j)
		
		
		
		
		

Table 4. The table showcases examples of the “FaceLifting” process. It is worth observing that the process of beautification prefers greenery, narrow roads and pavements

Q1 People's perceptions of beautified scenes

To ascertain whether FaceLifted scenes are perceived by individuals as they are supposed to, we run a crowd-sourcing experiment on Amazon Mechanical Turk. We randomly select 200 scenes, 100 beautiful and 100 ugly (taken at the bottom 10 and top 10 percentiles of the Trueskill's score distribution of Figure 2). Our framework then transforms each ugly scene into its beautified version, and each beautiful scene into its corresponding 'uglified'. These scenes are arranged into pairs, each of which contains the original scene and its beautified or uglified version. On Mechanical Turk, we only select verified masters for our crowd-sourcing workers (those with an approval rate above 90% during the past 30 days), pay them \$0.1 per task, and ask each of them to choose the beautiful scene for given pairs. We make sure to have at least 3 votes for each scene pair. Overall, our workers end up selecting the scenes that are actually beautiful 77.5% of the times, suggesting that FaceLifted scenes are correctly perceived most of the times.

Q2 Are beautified scenes great urban spaces?

To answer that question, we need to understand what makes a space great. After a careful review of the urban planning literature, we identify four factors [1, 11] (summarized in Table 5): great places mainly tend to be walkable, offer greenery, feel cozy, and be visually rich.

Metric	Description
Walkability	Walkable streets increase the social capital of a place, and they appeal to the exploring nature of the human psyche [11, 31, 39].
Green Spaces	The presence of greenery has repeatedly been found to impact people's well being [1]. Under certain conditions, it could also promote social interactions [32]. This suggest that not all greenery has to be considered in the same way though: dense forests or unkempt greens might well have a negative impact [16].
Landmarks	Losing a bearing in the city is not a very pleasant experience. Hence presence of recognisable and guiding landmarks influences the perception of an urban space [11, 25, 32].
Privacy-Openness	A sense of privacy (as opposed to a sense of openness) impacts a place's perception [11].
Visual Complexity	Visual complexity is a measure of how diverse a urban scene is in terms of design materials, textures, and objects [11]. We perceive complexity in an 'inverted-U' fashion, which means we prefer medium complexity over too little or too much for finding a place pleasant[43]

Table 5. Urban Design Metrics

To automatically extract visual cues related to these four factors, we select 500 ugly scenes and 500 beautiful ones at random, transform them into their opposite aesthetic qualities (i.e., ugly ones are beautified, and beautiful ones are 'uglified'), and compare which urban elements related to the four factors distinguish uglified scenes from beautified ones.

We extract labels from each of our 1,000 scenes using two image classifiers. First, using PlacesNet [46], we label each of our scenes according to a classification containing 205 labels (reflecting, for example, landmarks, natural elements), and retain the five labels with highest confidence scores for the scene. Second, using Segnet [4], we label each of our scenes according to a classification

containing 12 labels. Segnet is trained on dash-cam images, and classifies each scene pixel with one of these twelve labels: road, sky, trees, buildings, poles, signage, pedestrians, vehicles, bicycles, pavement, fences, and road markings.

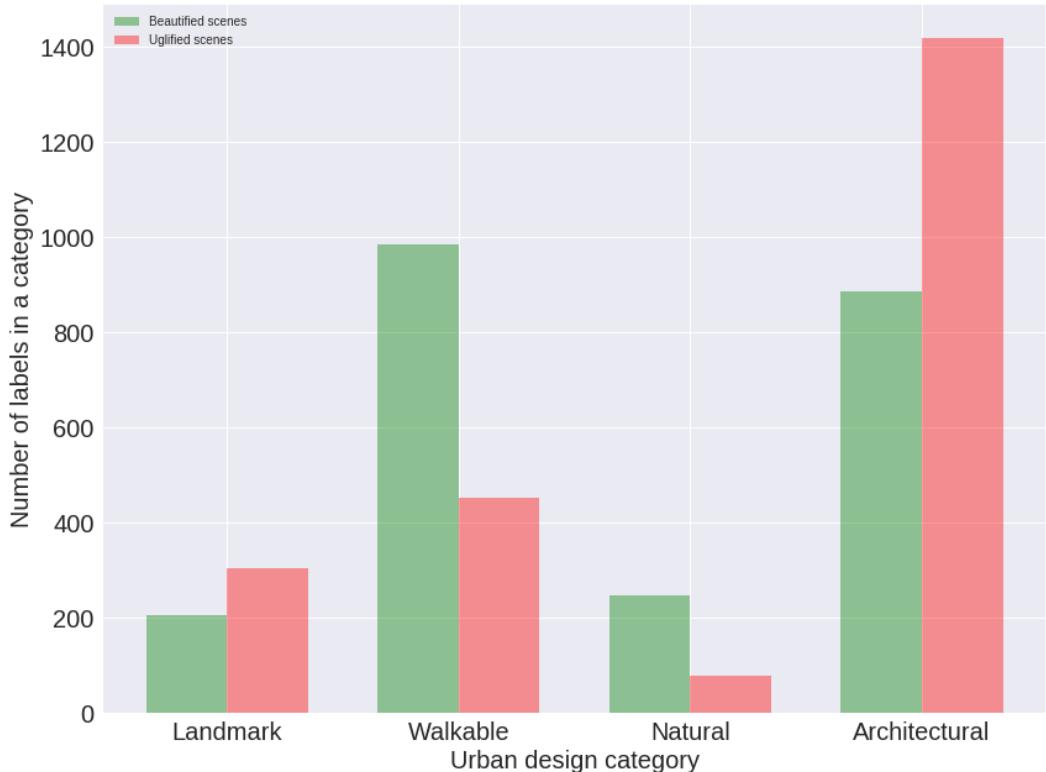


Fig. 6. Number of labels in specific urban design categories (on the *x*-axis) found in beautified scenes as opposed to those found in uglified scenes.

Having these two ways of labeling scenes, we can now test whether the expectations set by the literature describing metrics of great urban spaces (Table 5) are met in the FaceLifted scenes.

H1 Beautified scenes tend to be walkable. We manually select only the PlacesNet labels that are related to walkability. These labels include, for example, *abbey*, *plaza*, *courtyard*, *garden*, *picnic area*, and *park*. To test hypothesis *H1*, we count the number of walkability-related labels found in beautified scenes as opposed to those found in uglified scenes (Figure 6): the former contain twice as many walkability labels than the latter. We then determine which types of scenes are associated with beauty (Figure 7). Unsurprisingly, beautified scenes tend to show gardens, yards, and small paths. By contrast, uglified ones tend to show built environment features such as shop fronts and broad roads.

H2 Beautified scenes tend to offer green spaces. We manually select only the PlacesNet labels that are related to greenery. These labels include, for example, *fields*, *pasture*, *forest*, *ocean*, and *beach*. Then, in our 1,000 scenes, to test hypothesis *H2*, we count the number of nature-related labels found in beautified scenes as opposed to those found in uglified scenes (Figure 6): the former contain more

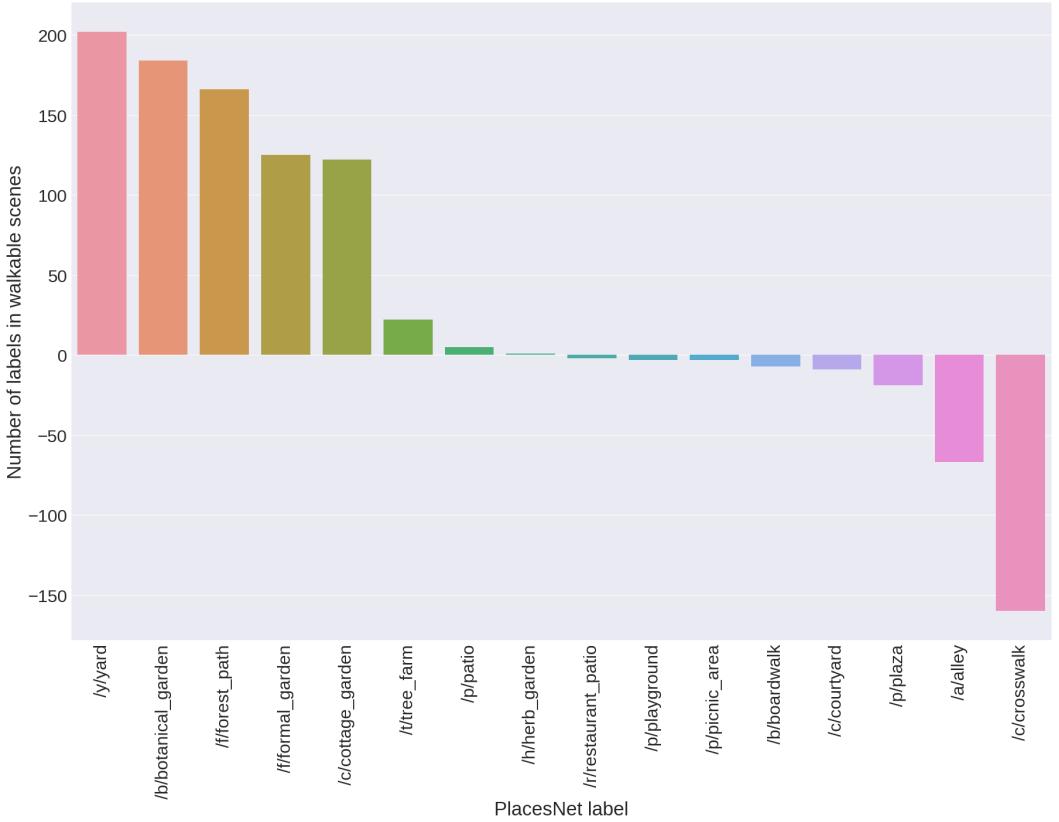


Fig. 7. Count of specific walkability-related labels (on the x -axis) found in beautified scenes minus the count of the same labels found in uglified scenes.

than twice as many nature-related labels than the latter. To test this hypothesis further, we compute the fraction of ‘tree’ pixels (using SegNet’s label ‘tree’) in beautified and uglified scenes, and find that beautification adds 32% of tree pixels, while uglification removes 17% of them.

H3 Beautified scenes tend to feel private and ‘cozy’. To test hypothesis H3, we count the fraction of pixels that Segnet labeled as ‘sky’ and show the results in a bin plot in Figure 8a: the x -axis has six bins (each of which represents a given range of sky fraction), and the y -axis shows the percentage of beautified vs. uglified scenes that fall into each bin. Beautified scenes tend to be cozier (lower sky presence) than the corresponding original scenes.

H4 Beautified scenes tend to be visually rich To quantify to which extent scenes are visually rich, we measure their visual complexity [11] as the amount of disorder in terms of distribution of (Segnet) urban elements in the scene:

$$H(X) = - \sum p(i) \log p(i) \quad (2)$$

where i is the i^{th} Segnet’s label. The total number of labels is twelve. The higher $H(X)$, the higher the scene’s entropy, that is, the higher the scene’s complexity. It has been proposed, that the perception of aesthetics or pleasantness follows an ‘inverted U’ shape[43]. To test hypothesis H4, we show

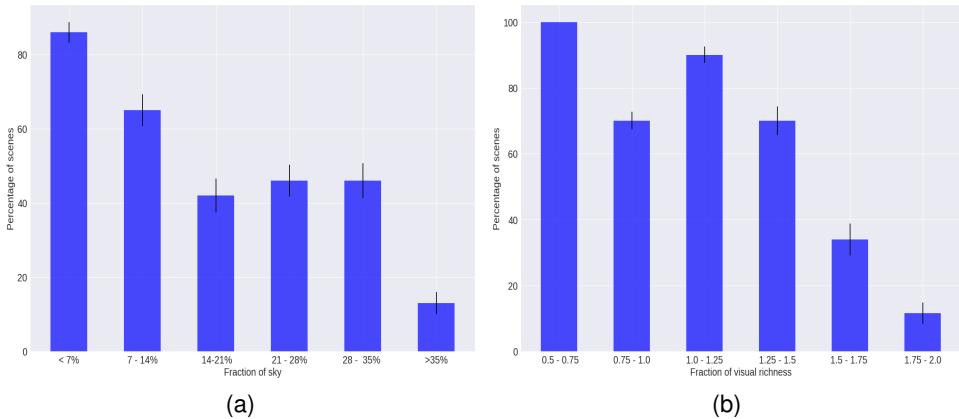


Fig. 8. The percentage of scenes (y -axis): (a) having an increasing presence of sky (on the x -axis); and (b) having an increasing level of visual richness (on the x -axis). The error bars represent standard errors obtained by random re-sampling of the data for 500 iterations.

Pair of urban elements	β_1	β_2	β_3	Error Rate (Percentage)
Buildings - Trees	-0.032	0.084	0.005	12.7
Sky - Buildings	-0.08	-0.11	0.064	14.4
Roads - Vehicles	-0.015	-0.05	0.023	40.6
Sky - Trees	0.03	0.11	-0.012	12.8
Roads - Trees	0.04	0.10	-0.031	13.5
Roads - Buildings	-0.05	-0.097	0.04	20.2

Table 6. Coefficients of logistic regressions run on one pair of predictors at the time.

the percentage of scenes that fall into a complexity bin (Figure 8b): beautified scenes are of low to medium complexity, while uglified ones are of high complexity.

Q3 Urban elements of beautified scenes

To determine which urban elements are the best predictors of urban beauty and the extent to which they are so, we run a logistic regression, and, to ease interpretation, we do so on one pair of predictors at the time:

$$Pr(\text{beautiful}) = \text{logit}^{-1}(\alpha + \beta_1 * V_1 + \beta_2 * V_2 + \beta_3 * V_1 \cdot V_2) \quad (3)$$

where V_1 is the fraction of the scene's pixels marked with one Segnet's label, say, "buildings" (over the total number of pixels), and V_2 is the fraction of the scene's pixels marked with another label, say, "trees". The result consists of three beta coefficients: β_1 reflects V_1 's contribution in predicting beauty, β_2 reflects V_2 's contribution, and β_3 is the interaction effect, that is, it reflects the contribution of the dependency of V_1 and V_2 in predicting beauty. We run logistic regressions on the five factors that have been found to be most predictive of urban beauty [1, 11, 32], and show the results in Table 6.

Since we are using logistic regressions, the quantitative interpretation of the beta coefficients is eased by the "divide by 4 rule" [44]: we can take β coefficients and "divide them by 4 to get an upper bound of the predictive difference corresponding to a unit difference" in beauty [44]. For example, take the results in the first row of Table 6. In the model $Pr(\text{beautiful}) = \text{logit}^{-1}(\alpha - 0.032 \cdot \text{buildings} + 0.084 \cdot \text{trees} + 0.005 \cdot \text{buildings} \cdot \text{trees})$, we can divide $-0.032/4$ to get -0.008 : a difference of 1 in the fraction of pixels being buildings corresponds to no more than a 0.8% negative difference

Use case	Definitely Not	Probably Not	Probably	Very Probably	Definitely
Decision Making	4.8%	9.5%	38%	28.6%	19%
Participatory Urban Planning	0%	4.8%	52.4%	23.8%	19%
Promote Green Cities	4.8%	0%	47.6%	19%	28.6%

Table 7. Urban experts polled about the extent to which an interactive map of “FaceLifted” scenes promotes: (a) decision making; (b) citizen participation in urban planning; and (c) promotion of green cities

in the probability of the scene being beautiful. In a similar way, a difference of 1 in the fraction of pixels being trees corresponds to no more than a 0.021% *positive* difference in the probability of the scene being beautiful. By considering the remaining results in Table 6, we find that, across all pairwise comparisons, trees is the most positive element associated with beauty, while roads and buildings are the most negative ones. Since these results go in the direction one would expect, one might conclude that the scenes beautified by our framework are in line with previous literature, adding further external validity to our work.

Q4 Do architects and urban planners find it useful?

To ascertain whether practitioners find FaceLift potentially useful, we built an interactive map of the city of Boston in which, for selected points, we showed pairs of urban scenes before/after beautification (Figure 9). We then sent that map along with a survey to 20 experts in architecture, urban planning, and data visualization around the world. Being experts in their respective fields, we wanted the survey takers to express a clear opinion about the utility of such a technology in their areas of practice, rather than express a non-committal neutral response. In accordance with this constraint, we designed the survey based on a non neutral response Likert scale, as explored in previous studies [5, 27]. The experts had to complete tasks in which they rated FaceLift based on how well it supports decision making, participatory urbanism, and promotion of green spaces among the general public. The results are show in Table 7 according to our experts, the tool can very probably supports decision making, probably support participatory urbanism, and definitely promote green spaces. These results are qualitatively supported by our experts’ comments, which include: “*The maps reveal patterns that might not otherwise be apparent*”, “*The tool helps focusing on parameters to identify beauty in the city while exploring it*”, and “*The metrics are nice. It made me think more about beautiful places needing a combination of criteria, rather than a high score on one or two dimensions. It made me realize that these criteria are probably spatially correlated*”.

5 CONCLUSION

FaceLift is a transparent framework that beautifies urban scenes. This translates into two main technical advancements. First, FaceLift is able to generate realistic beautified scenes based on existing approaches of Generative Adversarial Networks and deep convolutional networks. This is done so by implementing an activation maximization pipeline, that allows us to generate a maximized version of an originally un-attractive urban scene. Second, it augments the deep learning methodology with a module that offers explanations on what has been transformed and what are the predictors of the resulting beauty, in the language of practitioners. This makes such a methodology more understandable and adoptable as per our survey.

There are still important limitations though. One is data bias. The framework is as good as its training data, and more work has to go into collecting reliable ground truth of human perceptions. This data should ideally be stratified according to the people’s characteristics that impact their perceptions. The other main limitation is that generative models are hard to control, and more work has to go into offering principled ways of fine-tuning the generative process.

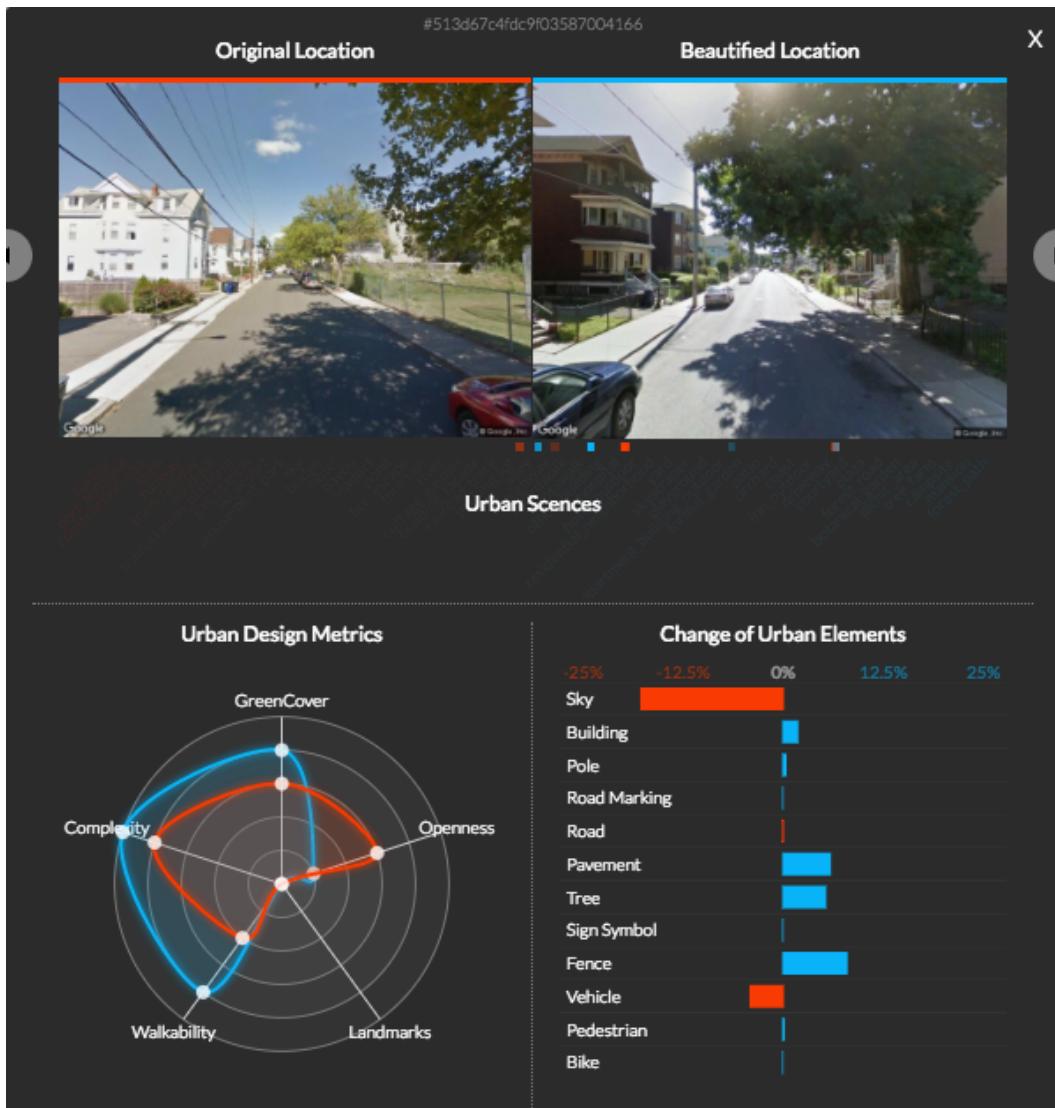


Fig. 9. Interactive map of FaceLifted scenes in Boston.

Despite these limitations, FaceLift has the potential to support urban interventions in scalable and replicable ways: it can be applied to the scale of an entire city, and that can be replicated in other cities. The advantage of shifting the focus of research away from predictive analytics towards urban interventions is that people could be part of discussions on works of architecture more than they are nowadays. To turn existing spaces into something more beautiful, that will still be the duty of architecture. Yet, with technologies similar to FaceLift more readily integrated in the architecture discussions, the complex job of recreating restorative spaces in an increasingly urbanized world will be greatly simplified. After all, “we delight in complexity to which genius have lent an appearance of

simplicity.” [7] In the context of future work, that genius is represented by the future technologies that we will contribute to build to deal with the complexity of our cities.

REFERENCES

- [1] Christopher Alexander, Sara Ishikawa, Murray Silverstein, Joaquim Romaguera i Ramió, Max Jacobson, and Ingrid Fiksdahl-King. 1977. *A Pattern Language: Towns, Buildings, Constructions*. Oxford University Press.
- [2] Hassan Abu Alhaija, Siva Karthik Mustikovela, Lars Mescheder, Andreas Geiger, and Carsten Rother. 2018. Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision* 126, 9 (2018), 961–972.
- [3] Sean M Arietta, Alexei A Efros, Ravi Ramamoorthi, and Maneesh Agrawala. 2014. City forensics: Using visual elements to predict non-visual city attributes. *IEEE transactions on visualization and computer graphics* 20, 12 (2014), 2624–2633.
- [4] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. 2015. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *arXiv preprint arXiv:1511.00561* (2015).
- [5] Aphrodite Baka, Lia Figgou, and Vasiliki Triga. 2012. ‘Neither agree, nor disagree’: a critical analysis of the middle answer category in Voting Advice Applications. *International Journal of Electronic Governance* 5, 3-4 (2012), 244–263.
- [6] Kylie Ball, Adrian Bauman, Eva Leslie, and Neville Owen. 2001. Perceived environmental aesthetics and convenience and company are associated with walking for exercise among Australian adults. *Preventive medicine* 33, 5 (2001), 434–440.
- [7] A. De Botton. 2008. *The Architecture of Happiness*. Knopf Doubleday Publishing Group.
- [8] Marco De Nadai, Radu Laurentiu Vieriu, Gloria Zen, Stefan Dragicevic, Nikhil Naik, Michele Caraviello, Cesar Augusto Hidalgo, Nicu Sebe, and Bruno Lepri. 2016. Are Safer Looking Neighborhoods More Lively?: A Multimodal Investigation into Urban Life. In *Proceedings of the ACM on Multimedia Conference (MM)*.
- [9] Alexey Dosovitskiy and Thomas Brox. 2016. Generating images with perceptual similarity metrics based on deep networks. In *Advances in Neural Information Processing Systems*. 658–666.
- [10] Abhimanyu Dubey, Nikhil Naik, Devi Parikh, Ramesh Raskar, and César A Hidalgo. 2016. Deep learning the city: Quantifying urban perception at a global scale. *arXiv preprint arXiv:1608.01769* (2016).
- [11] Reid Ewing and Otto Clemente. 2013. *Measuring urban design: Metrics for livable places*. Island Press.
- [12] Billie Giles-Corti, Melissa H Broomhall, Matthew Knuiman, Catherine Collins, Kate Douglas, Kevin Ng, Andrea Lange, and Robert J Donovan. 2005. Increasing walking: how important is distance to, attractiveness, and size of public open space? *American journal of preventive medicine* 28, 2 (2005), 169–176.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.
- [14] Ralf Herbrich, Tom Minka, and Thore Graepel. 2007. TrueSkillā‘: a Bayesian skill rating system. In *Advances in neural information processing systems*. 569–576.
- [15] Desislava Hristova, Luca M. Aiello, and Daniele Quercia. 2018. The New Urban Success: How Culture Pays. *Frontiers in Physics* 6 (2018), 27. DOI : <http://dx.doi.org/10.3389/fphy.2018.00027>
- [16] J. Jacobs. 1961. *The Death and Life of Great American Cities*. Random House.
- [17] Rachel Kaplan and Stephen Kaplan. 1989. *The experience of nature: A psychological perspective*. CUP Archive.
- [18] Stephen Kaplan, Rachel Kaplan, and John S Wendt. 1972. Rated preference and complexity for natural and urban visual material. *Perception & Psychophysics* 12, 4 (1972), 354–356.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [20] Stephen Law, Brooks Paige, and Chris Russell. 2018. Take a look around: using street view and satellite images to estimate house prices. *arXiv preprint arXiv:1807.07155* (2018).
- [21] Stephen Law, Chanuki Illushka Seresinhe, Yao Shen, and Mario Gutierrez-Roig. 2018. Street-Frontage-Net: urban image classification using deep convolutional neural networks. *International Journal of Geographical Information Science* 0, 0 (2018), 1–27. DOI : <http://dx.doi.org/10.1080/13658816.2018.1555832> arXiv:<https://doi.org/10.1080/13658816.2018.1555832>
- [22] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and others. 2017. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4681–4690.
- [23] Pall Jakob Lindal and Terry Hartig. 2012. Architectural variation, building height, and the restorative quality of urban residential streetscapes. *Journal of Environmental Psychology* (2012).

- [24] Pauline Luc, Camille Couprie, Soumith Chintala, and Jakob Verbeek. 2016. Semantic segmentation using adversarial networks. *arXiv preprint arXiv:1611.08408* (2016).
- [25] Kevin Lynch. 1960. *The image of the city*. Vol. 11.
- [26] Liqian Ma, Qianru Sun, Stamatios Georgoulis, Luc Van Gool, Bernt Schiele, and Mario Fritz. 2018. Disentangled person image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 99–108.
- [27] Guy Moors. 2008. Exploring the effect of a middle response category on response style in attitude measurement. *Quality and Quantity* 42, 6 (2008), 779–794.
- [28] Nikhil Naik, Scott Duke Kominers, Ramesh Raskar, Edward L Glaeser, and César A Hidalgo. 2017. Computer vision uncovers predictors of physical urban change. *Proceedings of the National Academy of Sciences* 114, 29 (2017), 7571–7576.
- [29] Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. 2016. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems*. 3387–3395.
- [30] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. 2016. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2536–2544.
- [31] Daniele Quercia, Luca Maria Aiello, Rossano Schifanella, and Adam Davies. 2015. The Digital Life of Walkable Streets. In *Proceedings of the 24th ACM Conference on World Wide Web (WWW)*. 875–884.
- [32] Daniele Quercia, Neil Keith O'Hare, and Henriette Cramer. 2014. Aesthetic capital: what makes London look beautiful, quiet, and happy?. In *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*. ACM, 945–955.
- [33] Daniele Quercia, Rossano Schifanella, and Luca Maria Aiello. 2014. The shortest path to happiness: Recommending beautiful, quiet, and happy routes in the city. In *Proceedings of the 25th ACM conference on Hypertext and social media*. ACM, 116–125.
- [34] Eulogio Real, Constantino Arce, and José Manuel Sabucedo. 2000. Classification of landscapes using quantitative and categorical data, and prediction of their scenic beauty in north-western Spain. *Journal of environmental psychology* 20, 4 (2000), 355–373.
- [35] Philip Salesses, Katja Schechtner, and César A Hidalgo. 2013. The collaborative image of the city: mapping the inequality of urban perception. *PLoS one* 8, 7 (2013), e68400.
- [36] Robert J. Sampson and Stephen W. Raudenbush. 2004. Seeing Disorder: Neighborhood Stigma and the Social Construction of Broken Windows. *Social Psychology Quarterly* 67, 4 (2004).
- [37] Chanuki Illushka Seresinhe, Tobias Preis, and Helen Susannah Moat. 2015. Quantifying the impact of scenic environments on health. *Scientific reports* 5 (2015), 16899.
- [38] Chanuki Illushka Seresinhe, Tobias Preis, and Helen Susannah Moat. 2017. Using deep learning to quantify the beauty of outdoor places. *Royal Society open science* 4, 7 (2017), 170170.
- [39] J. Speck. 2012. Walkable City: How Downtown Can Save America, One Step at a Time. In *Farrar, Straus and Giroux*.
- [40] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15, 1 (2014), 1929–1958.
- [41] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1–9.
- [42] Yaniv Taigman, Adam Polyak, and Lior Wolf. 2016. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200* (2016).
- [43] Roger S Ulrich. 1983. Aesthetic and affective response to natural environment. In *Behavior and the natural environment*. Springer, 85–125.
- [44] Brandon K Vaughn. 2008. Data analysis using regression and multilevel/hierarchical models, by Gelman, A., & Hill, J. *Journal of Educational Measurement* 45, 1 (2008), 94–97.
- [45] Xinchen Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee. 2015. Attribute2image: Conditional image generation from visual attributes. CoRR abs/1512.00570. (2015).
- [46] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. 2014. Learning deep features for scene recognition using places database. In *Advances in neural information processing systems*. 487–495.