# AgingMapGAN (AMGAN): High-Resolution Controllable Face Aging with Spatially-Aware Conditional GANs

Julien Despois, Frédéric Flament, Matthieu Perrot

L'Oréal Research & Innovation

**Abstract.** Existing approaches and datasets for face aging produce results skewed towards the mean, with individual variations and expression wrinkles often invisible or overlooked in favor of global patterns such as the fattening of the face. Moreover, they offer little to no control over the way the faces are aged and can difficultly be scaled to large images, thus preventing their usage in many real-world applications. To address these limitations, we present an approach to change the appearance of a high-resolution image using ethnicity-specific aging information and weak spatial supervision to guide the aging process. We demonstrate the advantage of our proposed method in terms of quality, control, and how it can be used on high-definition images while limiting the computational overhead.

**Keywords:** Conditional GANs, Face Aging, High-Resolution

## 1 Introduction



Fig. 1: High-resolution faces aged in a continuous manner with our approach

Face aging is an image synthesis task in which a reference image must be transformed to give the impression of a person of different age while preserving the identity and key facial features of the subject. When done correctly,

this process can be used in various domains, from predicting the future appearance of a missing person to entertainment and educational uses. We focus on achieving high-resolution face aging, as it is a required step towards capturing the fine details of aging (fine lines, pigmentation, etc.). In recent years, Generative Adversarial Networks [14] have allowed a learning-based approach for this task. The results, however, often lack in quality and only provide limited aging options. Popular models such as StarGAN [10] fail to produce convincing results without additional tweaks and modifications. This partially stems from the choice of reducing aging to one's real or apparent age [1]. Also, current approaches treat aging as a step-wise process, splitting age in bins (30-40, 40-50, 50+, etc.) [2,16,30,32,34].

In reality, aging is a continuous process that can take many forms depending on genetic factors such as facial features and ethnicity, as well as lifestyle choices (smoking, hydration, sun damage, etc.) or behavior. Notably, expression wrinkles are promoted by habitual facial expressions and can be prominent on the forehead, upper lip, or at the corner of the eyes (crow's feet). In addition, aging is subjective as it depends on the cultural background of the person assessing the age. These factors call for a more fine-grained approach to face aging.

In this paper, we aim to obtain high-resolution face aging results by creating a model capable of individually transforming the local aging signs. Our contributions are as follows:

- We show that a curated high-resolution dataset in association with a combination of novel and existing techniques produces detailed state-of-the-art aging results.
- We demonstrate how clinical aging signs and weak spatial supervision allows fine-grained control over the aging process of the different parts of the face.
- We introduce a patch-based approach to enable inference on high-resolution images while keeping the computational cost of training the model low.

## 2   Related Work

**Conditional Generative Adversarial Networks** Generative Adversarial Networks [14] leverage the principle of an adversarial loss to force samples generated by a generative model to be indistinguishable from real samples. This approach led to impressive results, especially in the domain of image generation. GANs can be extended to generate images based on one or several conditions. The resulting Conditional Generative Adversarial Networks are trained to generate images that satisfy both the realism and condition criteria.

**Unpaired Image-to-Image Translation** Conditional GANs are a powerful tool for image-to-image translation [18] tasks, where an input image is given to the model to synthesize a transformed image. StarGAN [10] introduced a way to

use an additional condition to specify the desired transformation to be applied. They propose to feed the input condition to the generator in the form of feature maps [10] concatenated to the input image, but new approaches use more complex mechanisms such as AdaIN [20] or its 2D extension SPADE [23] to give the generator the condition in a more optimal manner. Where previous techniques required pixel-aligned training images in the different domains, recent works such as CycleGAN [36] and StarGAN [10] introduced a cycle-consistency loss to enable unpaired training between discrete domains. This has been extended in [24] to allow translation between continuous domains.

**Face Aging**  To age a face from a single picture, traditional approaches use training data of either one [2, 16, 29, 32, 34, 35] or multiple images [27, 30] of the same person, along with the age of the person when the picture was taken. The use of longitudinal data, with multiple photos of the same person, offers less flexibility as it creates a heavy time-dependent constraint on the dataset collection.

The age is usually binned into discrete age groups (20-30, 30-40, 40-50, 50+, etc.) [2, 16, 32, 34], which frames the problem more simply, but limits the control over the aging process and doesn't allow the training to leverage the ordered nature of the groups. [29, 35] address this limitation by considering age as a continuous value. However, aging isn't objective because different skin types age differently, and different populations look for different signs of aging. Focusing on the apparent age as the guide for aging thus freezes the subjective point of view. Such approaches cannot be tailored to a population's perspective without requiring additional age estimation data from their point of view.

To improve the quality and level of details of the generated images, [34] use the attention mechanism from [24] in the generator. The generated samples are, however, low-definition images which are too coarse for real-world applications. Working at this scale hides some difficulties of generating realistic images, such as skin texture, fine lines, and the overall sharpness of the details.

## 3   Proposed Approach
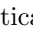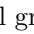
### 3.1   Problem Formulation

In this work, our goal is to use single unpaired images to train a model able to generate realistic high-definition ($1024 \times 1024$) aged faces, with continuous control over the fine-grained aging signs to create smooth transformations between the original and transformed images. This is a more intuitive approach, as aging is a continuous process and age group bins do not explicitly enforce a logical order.

We propose the use of ethnic-specific skin atlases [4–7, 13] to incorporate the ethnic dimension of clinical aging signs. These atlases define numerous clinical signs such as the wrinkles underneath the eye, the ptosis of the lower part of the face, the density of pigmentary spots on the cheeks, etc. Each sign is linked to a

specific zone on the face and scored on a scale that depends on ethnicity. Using these labels in addition to the age make for a more complete representation of aging, and allows transforming images with various combination of clinical signs and scores.

The aging target is passed to the network in the form of an aging map (Fig. 2). To do so, we compute facial landmarks and define the relevant zone for each aging sign. Each zone (e.g. forehead) is then filled with the score value of the corresponding sign (e.g. forehead wrinkles). We use the apparent age to fill in the blanks where the clinical signs are not defined. Finally, we coarsely mask the background of the image.

Treating the whole image at once would be ideal, but training a model with $1024 \times 1024$ images requires large computational resources. Our approach allows us to train the model by patch, using only part of the image during training, and the corresponding part of the aging map. Patch-based training reduces the context (i.e. global information) for the task but also reduces the computational resources required to process high-resolution images in large batches, as recommended in [8]. We leverage this to use a large batch size on small patches of $128 \times 128$, $256 \times 256$ or $512 \times 512$ pixels.

The major drawback of the patch-based training is that small patches can look similar (e.g. forehead and cheek) yet must be aged differently (e.g. respectively horizontal and vertical wrinkles). To avoid "mean" wrinkles on these ambiguous zones, we give the generator two patches coming respectively from a horizontal and a vertical gradient location map (■, ■). This allows the model to know the position of the patch in order to differentiate between potentially ambiguous zones.
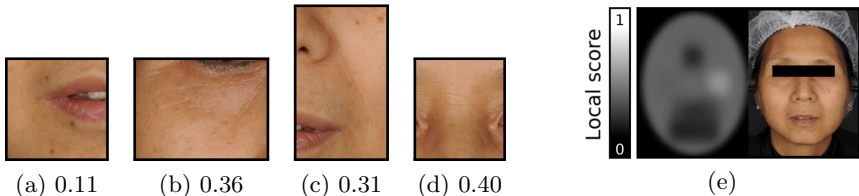


Fig. 2: Aging sign zones (a-d) and their associated scores used to construct the aging map (e). The brightness of each pixel represents the normalized score of the localized clinical sign (wrinkles at the corner of the lips (a), underneath the eye wrinkles (b), nasolabial fold wrinkles (c), inter-ocular wrinkles (d), etc.) or the age where no sign is defined

## 3.2    AMGAN - Network Architectures

We base our training process on the StarGAN [10] framework. Our generator is a fully convolutional encoder-decoder derived from [11] with SPADE [23] residual

blocks in the decoder to incorporate the aging and location maps. This allows the model to leverage the spatial information present in the aging map, and use it at multiple scales in the decoder. To avoid learning unnecessary details, we use the attention mechanism from [24] to force the generator to transform the image only where needed. The discriminator is a modified version of [10], and produces the outputs for the WGAN [3] objective (given for an image $i$ and aging map $a$ in Equation 1), the estimation of the coordinates of the patch, and the low-resolution estimation of the aging map. Fig. 3 and Fig. 4 present the patch-based training workflow.

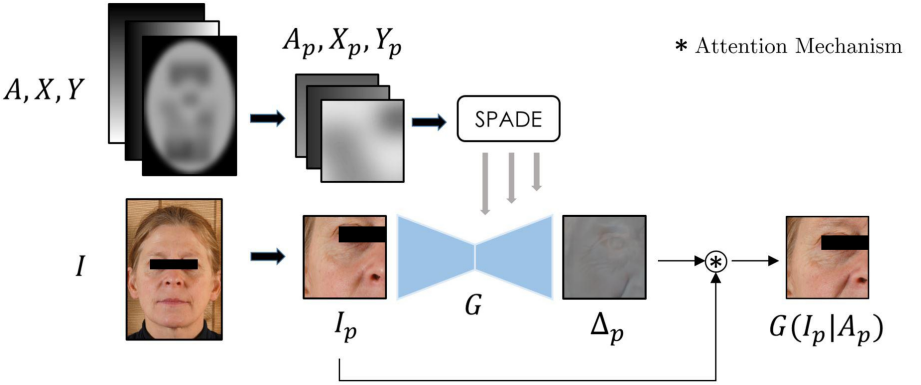$$\mathcal{L}_{WGAN} = \mathbb{E}_i[D(i)] - \mathbb{E}_{i,a}[D(G(i,a))] \tag{1}$$



Fig. 3: Generator of our proposed patch-based training. We begin by cropping a patch from the image $I$, aging map $A$, and location maps $X$ and $Y$. The generator transforms the image patch $I_p$ according to the map and location

### 3.3  Aging Maps

To avoid penalizing the model for failing to place the bounding boxes with pixel-perfect precision, we blur the aging maps to smooth the edges and compute the discriminator regression loss on downsampled $10 \times 10$ maps. This formulation allows packing the information in a more compact and meaningful way than as individual uniform feature maps [10, 30, 34, 35]. Our approach only requires multiple feature maps when there are large overlaps between signs (e.g. forehead pigmentation and forehead wrinkles). Considering an image patch $i$ and aging map patch $a$, the loss is given in Equation 2.

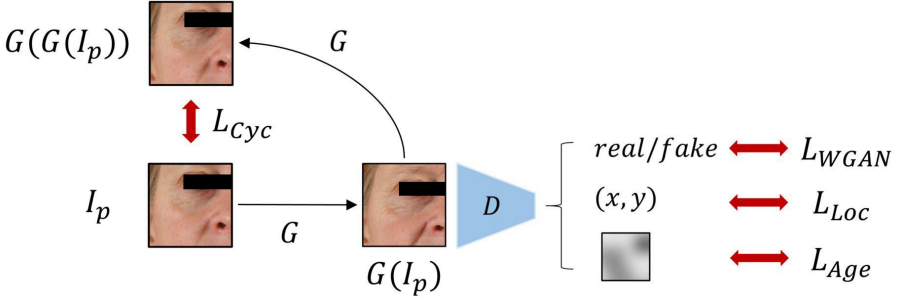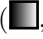$$\mathcal{L}_{Age} = \mathbb{E}_i[\|a - D_{Age}(G(i,a))\|_2] \tag{2}$$

Fig. 4: The discriminator produces the real/fake output, the estimated location of the patch, and the estimated local aging map. The outputs are respectively penalized with the WGAN objective, location, and aging map loss functions. The cycle consistency loss ensures the transformation preserves the key features of the original image

### 3.4   Location Maps

As in [22], we use two orthogonal gradients as location maps ( , ) to help the generator apply relevant aging transformations to a given patch. The $(x, y)$ coordinates of the patch could be given to the generator as two numbers instead of linear gradients maps, but doing so would prevent the use of the model on the full-scale image as it would break its fully-convolutional nature. Considering an image patch $i$ and aging map patch $a$ located at coordinates $(x, y)$, the loss is given in Equation 3.

$$\mathcal{L}_{Loc} = \mathbb{E}_i[\|(x, y) - D_{Loc}(G(i, a))\|_2] \tag{3}$$

### 3.5   Training

The models are trained with the Adam [21] optimizer with $\beta_1 = 0$, $\beta_2 = 0.99$ and learning rates of $7 \times 10^{-5}$ for G and $2 \times 10^{-4}$ for D. Following the two time-scale update rule [17], both models are updated at each step. Additionally, learning rates for both G and D are linearly decayed to zero over the course of the training. To enforce cycle-consistency, we use the perceptual loss [33] with $\lambda_{Cyc} = 100$. For the regression tasks, we use $\lambda_{Loc} = 50$ to predict the (x,y) coordinates of the patch and $\lambda_{Age} = 100$ to estimate the downsampled aging map. The discriminator is penalized with the original gradient penalty presented in [15] with $\lambda_{GP} = 10$. Our complete loss objective function is given in Equation 4.

$$\mathcal{L} = \mathcal{L}_{WGAN} + \lambda_{Cyc}\mathcal{L}_{Cyc} + \lambda_{Age}\mathcal{L}_{Age} + \lambda_{Loc}\mathcal{L}_{Loc} + \lambda_{GP}\mathcal{L}_{GP} \tag{4}$$

### 3.6 Inference

For inference, we use exponential moving average [31] over G's parameters. The trained generator can be used directly on the $1024 \times 1024$ image no matter the size of the patch used during training thanks to the fully convolutional nature of the network and the use of continuous 2D aging maps. We can either create a target aging map manually or use the face landmarks and target scores to build one.

## 4 Experiments

### 4.1 Experimental Setting

Most face aging datasets [9, 25, 26] suffer from a lack of diversity in terms of ethnicity [19], and focus on low-resolution images (up to $250 \times 250$ pixels). This isn't sufficient to capture fine details related to skin aging. Moreover, they often fail to normalize the pose and expression of the faces (smiling, frowning, raised eyebrows), which results in accentuated wrinkles unrelated to aging (mostly nasolabial wrinkles, crow's feet wrinkles, forehead wrinkles and wrinkles underneath the eye). Finally, the lack of fine-grained information on the aging signs causes other approaches to capture unwanted correlated features such as the fattening of the face, as observed in datasets such as IMDB-Wiki [26]. These effects can be observed in Fig. 5.

To address these issues, we tested our models on two curated high-resolution datasets, using generated aging maps or uniform aging maps to highlight the rejuvenation/aging.

### 4.2 Flicker Faces High-Quality Dataset (FFHQ)

To tackle the high-resolution aging problem, we have tested our approach on the FFHQ dataset [20]. To minimize the issues in lighting, pose, and facial expressions, we applied simple heuristics to select a subset of the dataset of better quality. To do so, we extracted facial landmarks from all faces and used them to remove all images where the head was too heavily tilted left, right, up, or down. In addition, we removed all images with an open mouth to limit artificial nasolabial fold and underneath the eye wrinkles. Finally, we used a HOG [12] feature descriptor to remove images with hair covering the face. This selection brings down the dataset from 70k+ to 10k+ images. Due to the extreme diversity of the FFHQ dataset, the remaining images are still far from being perfect, especially in terms of lighting color, direction, and exposure.

To obtain the scores of the individual aging signs on these images, we used aging sign estimation models based on the ResNet [28] architecture that we trained on the dataset described in Section 4.3. Finally, we generated the ground truth aging maps using the landmarks as a basis for the coarse bounding-boxes. We trained our model on $256 \times 256$ patches randomly selected on the $1024 \times 1024$ face.

|  | Original | 20-30 | 30-40 | 40-50 | 50-60 | 60+ |

(a) Wang et al. [30]

(b) Heljakka et al. [16]

(c) Song et al. [27]

(d) Antipov et al. [2]
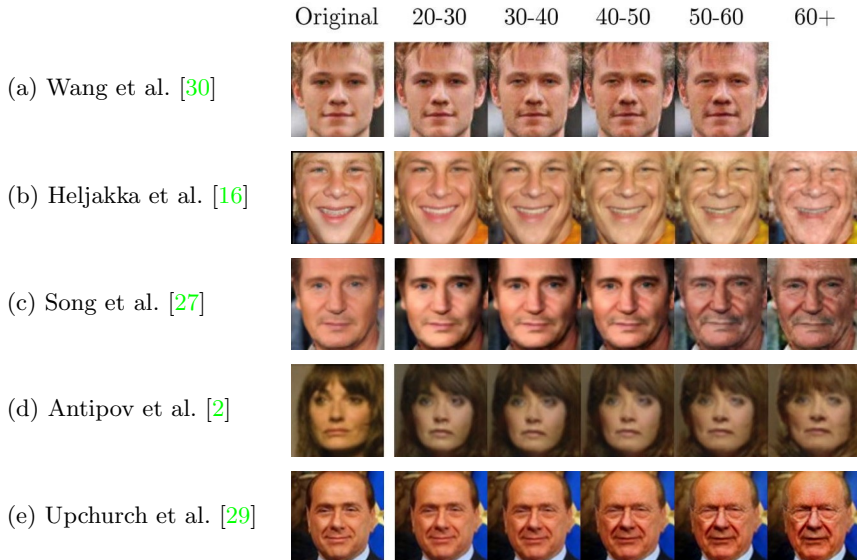
(e) Upchurch et al. [29]

Fig. 5: Previous approaches operate on low-resolution images and suffer from a lack of wrinkles dynamic range, especially for expression wrinkles (a). They are also prone to color shifts and artifacts (b, c, d, e), as well as unwanted correlated features such as the fattening of the face (d), or the addition of glasses (e)

### 4.3   High-Quality Standardized Dataset

To obtain better performance, we have collected a dataset of 6000 high-resolution ($3000 \times 3000$) images of faces, centered and aligned, spanning most ages (18-80), genders, and ethnicities (African, Caucasian, Chinese, Japanese and Indian). The images were labeled using ethnicity-specific clinical aging sign atlases [4–7, 13] and scored on signs covering most of the face (apparent age, forehead wrinkles, nasolabial fold, underneath the eye wrinkles, upper lip wrinkles, wrinkles at the corner of the lips and ptosis of the lower part of the face).

## 5   Results

### 5.1   FFHQ

Despite the complexity of the dataset, and without ground truth age values, our patch-based model is able to transform the individual wrinkles on the face in a continuous manner.

Fig. 6 displays how the model was able to transform the different wrinkles despite the complexity of the patch-based training, the large variation in lighting in the dataset, and the unbalance between grades of clinical signs/age, with a vast majority of young subjects with few wrinkles. Fig. 7 highlights the control

Fig. 6: Rejuvenation (*left*), original (*center*) and aging (*right*) for faces of different age and ethnicity from FFHQ dataset using our approach



Fig. 7: Where no sign is defined, we fill the map with the age. This helps the model learn global features like the greying of the hair (*left*). Using individual clinical signs in an aging map allows us to age all signs but keep the appearance of the hair intact (*right*)

we have over the individual signs, allowing aging the face in a controllable way that wouldn't be possible with the only label of the age.

## 5.2 High-Quality Standardized Dataset

On our standardized images, and with better coverage across ethnicity and aging signs, our model demonstrates state-of-the-art performance (Fig. 1, Fig. 8), with a high level of detail, realism, and no visible artifacts, color shifts or unwanted correlated features as seen in previous works (Fig. 5). The aging process is successful along the continuous spectrum of age maps, allowing realistic images to be generated for a diverse set of sign severity values (Fig. 9). More examples as well as HD videos are available in the supplementary materials and at https://despoisj.github.io/AgingMapGAN/.

Fig. 8: Faces aged in a continuous manner. No zone is left unchanged, even the forehead or the sagging of lower part of the face. The complementary age information used to fill the gap can be seen on the thinning or greying of the eyebrows. Note: zooming is recommended to see the fine details of the figure



(a)          (b)          (c)          (d)          (e)          (f)

Fig. 9: Same face (a) aged with different aging maps. (b) rejuvenates all signs except for the nasolabial, corner of the lips and underneath the eyes wrinkles on the right part of the face. (c) only ages the bottom of the face and (d) only the top. (e) only ages the wrinkles underneath the eye. (f) ages the face in an asymmetric fashion, namely the right wrinkles underneath the eyes and the left nasolabial fold

## 5.3    Evaluation Metrics

To be considered successful, the task of face aging requires three criteria to be met: the image must be realistic, the identity of the subject must be preserved, and the face must be aged. These are respectively enforced during training thanks to the WGAN objective function, cycle-consistency loss, and aging map estimation loss. By nature, one single metric couldn't ensure that all criteria are met. For instance, the model could leave the input image without altering it, and still succeed in realism and identity. Contrarily, the model could succeed in aging but fail realism and/or identity. If one model isn't superior to another on every metric, we need to choose a trade-off.

Our experiments on FFHQ and our high-quality standardized dataset never displayed any issue in the preservation of the subject identity. This is a consequence of the use of the attention mechanism that allows the generator to preserve the key facial features of the face. As a result, we chose to focus on the realism and aging criteria for our quantitative evaluation. Because our approach focuses on aging as a combination of aging signs instead of relying solely on age, we don't use the accuracy of the target age as a metric. Instead, we use the Fréchet Inception Distance (FID) [17] to assess the realism of our images, and the Mean Average Error (MAE) for the accuracy of the target aging signs.

To do so, we use half of our dataset as a reference for real images, and the rest as the images to be transformed by our model. The aging maps used to transform these images are chosen randomly from the ground truth labels to ensure a distribution of generated images that follows the original dataset. We estimate the value of individual scores on all generated images using dedicated aging sign estimation models based on the ResNet [28] architecture. As a reference for the FID scores, we compute the FID between both halves of the real image dataset. Note that the size of our dataset prevents us from computing the FID on the recommended 50k+ [17,20], thus leading to the overestimation of the value. This can be seen when computing the FID between real images only, giving a baseline FID of 49.0. The results are presented in Table 1.

Table 1: Fréchet Inception Distance and Mean Average Error for our different models

| Method | Patch Size | FID↓ | MAE↓ |
|:------:|:----------:|:----:|:----:|
| Real Images | - | 49.0 | - |
| AMGAN (Ours) | $512 \times 512$ | **110.1** | **0.14** |
| AMGAN (Ours) | $256 \times 256$ | 110.7 | **0.14** |
| w/o Aging Maps | $256 \times 256$ | 141.6 | 0.17 |
| AMGAN (Ours) | $128 \times 128$ | 112.9 | 0.17 |
| w/o Location Maps | $128 \times 128$ | 140.0 | 0.20 |

## 5.4    Comparison Between Age and Clinical Signs

When trained without clinical signs, using only the age to create a uniform aging map, the model still gives convincing results, with low FID and estimated age MAE. (Table 2).

By comparing the results with the full aging map approach, however, it appears that some wrinkles don't exhibit their full range of dynamics. This is due to the fact that not all aging signs need to be maximized in order to reach the limit age of the dataset. In fact, the 150 oldest individuals of our standardized dataset (65 to 80 years old) display a median standard deviation of their normalized aging signs of 0.18, highlighting the many possible combinations of aging signs in old people (Supplementary Materials, Fig. 1). For example, signs such as the forehead wrinkles are highly dependant on the facial expressions of the subject and are integral parts of the aging process. This an issue for the age-only model because it only offers one way to age a face.

To the contrary, the faces aged with the aging map offer much more control over the aging process. By controlling each individual sign of aging, we can choose whether to apply these expression wrinkles or not. A natural extension of this benefit is the pigmentation of the skin, which is viewed in some Asian countries as a sign of aging. An age-based model cannot produce aging for these countries without having to re-estimate the age from the local perspective. This doesn't scale, unlike our approach which, once trained with every relevant aging sign, can offer a face aging experience customized to the point of view of different countries, all in a single model and without additional labels.

Table 2: Fréchet Inception Distance and Mean Average Error for our model with clinical signs, and with age only

| Method | Patch Size | Control | FID | MAE |
|---|---|---|---|---|
| AMGAN (Ours) | $256 \times 256$ | ✓ | 110.7 | 0.143 |
| w/o Clinical Signs | $256 \times 256$ | ✗ | 101.3 | 0.116 |

## 5.5    Ablation Study

**Effect of Patch Size**  When training the model for a given target image resolution ($1024 \times 1024$ pixels in our experiments), we can choose the size of the patch used for the training. The bigger the patch, the more context the model will have to perform the aging task. For the same computation power, however, larger patches cause the batch size to be smaller, which hinders the training [8]. We conducted experiment using patches of $128 \times 128$, $256 \times 256$ and $512 \times 512$ pixels. Fig. 10 shows that all patch sizes manage to age the high-resolution face but to various degrees of realism. The smallest patch size suffers most from the

lack of context and produces results that are inferior to the other two, with visible texture artifacts. The $256 \times 256$ patch gives convincing results, with minor imperfection only visible when compared to the $512 \times 512$ patch. These results suggest that we could apply this technique to larger resolutions, such as with patches of $512 \times 512$ on $2048 \times 2048$ images.



Fig. 10: Results of rejuvenation (*left*) and aging (*right*) for different patch sizes on a $1024 \times 1024$ image

**Location Maps** To see the contribution of the location maps, we compared our model trained with and without them. As expected, the effect of the location maps is more prominent on small patch sizes, where the ambiguity is high. Fig. 11 shows how on small patch sizes and in the absence of location information, the model is unable to differentiate similar patches from different parts of the face. It is, therefore, unable to add wrinkles that are coherent with the location, and generates generic diagonal ripples. This effect is less present on larger patch sizes because the location of the patch is less ambiguous.

**Spatialization of Information** We compare our proposed aging maps against the baseline method of formatting conditions, namely to give all sign scores as individual uniform feature maps. Since not every sign is present in the patch, especially when the patch size is small, most of the processed information is of no use to the model. The aging maps represent a simple way of only giving the model the labels present in the patch, in addition to their spatial extent and location. Fig. 12 highlights the effect of the aging map. On small patches ($128 \times 128$, $256 \times 256$ pixels), the model struggles to create realistic results. The aging map helps reduce the complexity of the problem.

Fig. 11: Face aged with smallest patch size without (*left*) and with (*right*) location maps, along with the difference with the original image. The location maps eliminate the presence of diagonal texture artifacts, especially on the forehead where they allow horizontal wrinkle to appear



Fig. 12: Face aged with medium patch size with individual uniform condition feature maps (*left*) and proposed aging maps (*right*), along with the difference with the original image. The aging maps help make the training more efficient thanks to denser spatialized information, and produce more realistic aging. The difference highlights the small unrealistic wrinkles for the baseline technique

## 6   Conclusion

In this paper, we presented the use of clinical signs to create aging maps for face aging. Thanks to this technique, we demonstrated state-of-the-art results on high-resolution images with complete control over the aging process. Our patch-based approach allows conditional generative adversarial networks to be trained on large images while keeping a large batch size. This technique is applicable to various problems and can be used to tackle high-resolution problems with limited computational resources. In the future, the use of longitudinal data following the same person over time would allow a better understanding of the evolution of the aging signs on an individual basis, and therefore better personalizing of face aging factoring lifestyle, environmental and behavioral components.

# References

1. Agustsson, E., Timofte, R., Escalera, S., Baro, X., Guyon, I., Rothe, R.: Apparent and real age estimation in still images with deep residual regressors on appareal database. In: 2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017). pp. 87–94. IEEE (2017)
2. Antipov, G., Baccouche, M., Dugelay, J.L.: Face aging with conditional generative adversarial networks. In: 2017 IEEE international conference on image processing (ICIP). pp. 2089–2093. IEEE (2017)
3. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein gan. arXiv preprint arXiv:1701.07875 (2017)
4. Bazin, R., Doublet, E.: Skin aging atlas. volume 1. caucasian type. MEDCOM publishing (2007)
5. Bazin, R., Flament, F.: Skin aging atlas. volume 2, asian type (2010)
6. Bazin, R., Flament, F., Giron, F.: Skin aging atlas. volume 3. afro-american type. Paris: Medcom (2012)
7. Bazin, R., Flament, F., Rubert, V.: Skin aging atlas. volume 4, indian type (2015)
8. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. arXiv preprint arXiv:1809.11096 (2018)
9. Chen, B.C., Chen, C.S., Hsu, W.H.: Cross-age reference coding for age-invariant face recognition and retrieval. In: European conference on computer vision. pp. 768–783. Springer (2014)
10. Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8789–8797 (2018)
11. Choi, Y., Uh, Y., Yoo, J., Ha, J.W.: Stargan v2: Diverse image synthesis for multiple domains. arXiv preprint arXiv:1912.01865 (2019)
12. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: 2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05). vol. 1, pp. 886–893. IEEE (2005)
13. Flament, F., Bazin, R., Qiu, H.: Skin aging atlas. volume 5, photo-aging face & body (2017)
14. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in neural information processing systems. pp. 2672–2680 (2014)
15. Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C.: Improved training of wasserstein gans. In: Advances in neural information processing systems. pp. 5767–5777 (2017)
16. Heljakka, A., Solin, A., Kannala, J.: Recursive chaining of reversible image-to-image translators for face aging. In: International Conference on Advanced Concepts for Intelligent Vision Systems. pp. 309–320. Springer (2018)
17. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: Advances in neural information processing systems. pp. 6626–6637 (2017)
18. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1125–1134 (2017)
19. Kärkkäinen, K., Joo, J.: Fairface: Face attribute dataset for balanced race, gender, and age. arXiv preprint arXiv:1908.04913 (2019)

20. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4401–4410 (2019)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
22. Liu, R., Lehman, J., Molino, P., Such, F.P., Frank, E., Sergeev, A., Yosinski, J.: An intriguing failing of convolutional neural networks and the coordconv solution. In: Advances in Neural Information Processing Systems. pp. 9605–9616 (2018)
23. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2337–2346 (2019)
24. Pumarola, A., Agudo, A., Martinez, A.M., Sanfeliu, A., Moreno-Noguer, F.: Ganimation: Anatomically-aware facial animation from a single image. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 818–833 (2018)
25. Ricanek, K., Tesafaye, T.: Morph: A longitudinal image database of normal adult age-progression. In: 7th International Conference on Automatic Face and Gesture Recognition (FGR06). pp. 341–345. IEEE (2006)
26. Rothe, R., Timofte, R., Van Gool, L.: Dex: Deep expectation of apparent age from a single image. In: Proceedings of the IEEE international conference on computer vision workshops. pp. 10–15 (2015)
27. Song, J., Zhang, J., Gao, L., Liu, X., Shen, H.T.: Dual conditional gans for face aging and rejuvenation. In: IJCAI. pp. 899–905 (2018)
28. Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A.: Inception-v4, inception-resnet and the impact of residual connections on learning. In: Thirty-first AAAI conference on artificial intelligence (2017)
29. Upchurch, P., Gardner, J., Pleiss, G., Pless, R., Snavely, N., Bala, K., Weinberger, K.: Deep feature interpolation for image content changes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7064–7073 (2017)
30. Wang, Z., Tang, X., Luo, W., Gao, S.: Face aging with identity-preserved conditional generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7939–7947 (2018)
31. Yazici, Y., Foo, C.S., Winkler, S., Yap, K.H., Piliouras, G., Chandrasekhar, V.: The unusual effectiveness of averaging in gan training. arXiv preprint arXiv:1806.04498 (2018)
32. Zeng, H., Lai, H., Yin, J.: Controllable face aging. arXiv preprint arXiv:1912.09694 (2019)
33. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 586–595 (2018)
34. Zhu, H., Huang, Z., Shan, H., Zhang, J.: Look globally, age locally: Face aging with an attention mechanism. arXiv preprint arXiv:1910.12771 (2019)
35. Zhu, H., Zhou, Q., Zhang, J., Wang, J.Z.: Facial aging and rejuvenation by conditional multi-adversarial autoencoder with ordinal regression. arXiv preprint arXiv:1804.02740 (2018)
36. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Proceedings of the IEEE international conference on computer vision. pp. 2223–2232 (2017)