

Masked Linear Regression for Learning Local Receptive Fields for Facial Expression Synthesis

Nazar Khan Arbish Akram Arif Mahmood Sania Ashraf
Kashif Murtaza

{nazarkhan, phdcsf18m002, sashraf, kashifmurtaza}@pucit.edu.pk
arif.mahmood@itu.edu.pk
<http://faculty.pucit.edu.pk/nazarkhan/>

Invited Talk at ITU, Lahore, Pakistan
15th October 2019

We present an FES method with the following contributions:

- 1 Convex optimization with closed-form solution of global minimum in a single iteration.
- 2 Extremely low spatial and computational complexity.
- 3 Trainable on very small datasets.
- 4 Intuitive interpretation of learned parameters can be exploited to improve results.
- 5 Good generalization over different types of images that state-of-the-art GANs find very challenging to synthesize.

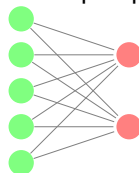
Related Work

- Basis learning (Blanz, Vetter, et al. 1999)
- Active appearance models (Cootes, Edwards, Taylor, et al. 2001)
- Deep belief nets (Susskind et al. 2008)
- Kernel regression (Huang and De la Torre 2010)
- GANs for image-to-image translation
 - Pix2Pix (Isola et al. 2017)
 - CycleGAN (Zhu et al. 2017)
 - StarGAN (Choi et al. 2018)
 - GANimation (Pumarola et al. 2019)

Regression

- Let $\mathbf{x} \in \mathbb{R}^D$ be a vectorized input image.
- Let $\mathbf{y} \in \mathbb{R}^K$ be a vectorized output image.
- Standard linear regression (ℓ_2) models output as $\mathbf{y} = W\mathbf{x}$ where $W \in \mathbb{R}^{K \times D}$ is a transformation matrix.
- This model corresponds to **global receptive fields**.
- Each output pixel is produced by *looking at* all input pixels.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_K \end{bmatrix} = \begin{bmatrix} w_{11} & \dots & w_{1D} \\ w_{21} & \dots & w_{2D} \\ \vdots & & \vdots \\ w_{K1} & \dots & w_{KD} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$



Global

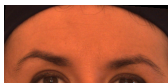
ℓ_2 -regression – error formulation

$$E^{\text{RR}}(W) = \frac{1}{2} \|WX^T - T^T\|_F^2 + \frac{\lambda_2}{2} \|W\|_F^2 \quad (1)$$

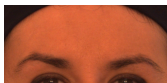
- $X \in \mathbb{R}^{N \times D}$ and $T \in \mathbb{R}^{N \times K}$ are the design matrices of vectorized input and target images respectively.
- Regularization parameter $\lambda_2 > 0$ controls over-fitting and $\|\cdot\|_F^2$ is the squared Frobenious norm of a matrix.
- This is a quadratic optimization problem with a global minimizer obtained in closed-form as

$$W^{\text{RR}} = ((X^T X + \lambda_2 I)^{-1} X^T T)^T \quad (2)$$

Do all pixels determine expression?



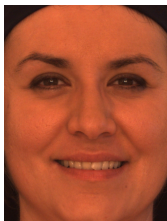
Expression=?



Expression=?



Neutral



Happy

- Is there any benefit of *looking at* forehead pixels to generate smiling lips?
- Happy lips can be generated by looking at and transforming lips.
- Happy eyes can be generated by looking at and transforming eyes.
- So why carry so many parameters in W ?

Expressions are local

- Transformation from one facial expression to another depends more on local information and less on global information.
- Facial expressions often constitute **sparsely distributed** and **locally correlated** changes.



Neutral

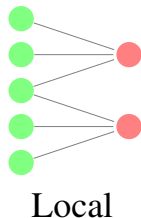


Happy

Masked Regression

- We propose a **Masked Regression (MR)**¹ model
 $\mathbf{y} = (W \circ M)\mathbf{x}$ where binary matrix M contains 1s only for locations that need to be looked at.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_K \end{bmatrix} = \begin{bmatrix} w_{11}m_{11} & \dots & w_{1D}m_{1D} \\ w_{21}m_{21} & \dots & w_{2D}m_{2D} \\ \vdots & \vdots & \vdots \\ w_{K1}m_{K1} & \dots & w_{KD}m_{KD} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$



- If $m_{ij} = 0$, then output pixel y_i is produced without looking at input pixel x_j .

¹N. Khan et al. "Masked Linear Regression for Learning Local Receptive Fields for Facial Expression Synthesis". In: *International Journal of Computer Vision (IJCV)* (2019).

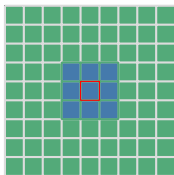
Linear Regression:

$$y_i = \sum_{j=1}^D w_{ij} x_j \quad (3)$$

Masked Regression:

$$y_i = \sum_{m_{ij}=1} w_{ij} x_j \quad (4)$$

If y_i is formed by looking at a 3×3 region in the input image, then the summation in MR is only over 9 pixels, irrespective of image size.

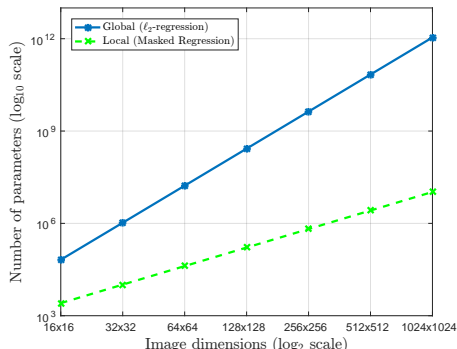


This corresponds to having **local receptive fields**.

		Input pixel index j in row-major order																								
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
Output pixel index i in row-major order	1	1	1				1	1																		
	2	1	1	1			1	1	1																	
	3		1	1	1			1	1	1																
	4			1	1	1			1	1	1															
	5				1	1				1	1															
	6	1	1				1	1				1	1													
	7	1	1	1			1	1	1			1	1	1												
	8		1	1	1			1	1	1			1	1	1											
	9			1	1	1			1	1	1			1	1	1										
	\vdots																									
	19													1	1	1			1	1	1			1	1	1
	20														1	1				1	1				1	1
	21															1	1				1	1				
	22															1	1	1			1	1	1			
	23																1	1	1			1	1	1		
	24																	1	1	1			1	1	1	
	25																		1	1				1	1	

Figure: Mask M corresponding to input image of size 5×5 , output image of size 5×5 and receptive fields of size 3×3 . For clarity, entries equal to 0 are left blank. If the entry at row i and column j is 1, then output pixel i has input pixel j in its receptive field.

Benefit of using mask



- Local receptive fields remain practical for larger image sizes.
- Regression with global receptive fields becomes impractical even for image sizes as small as 128×128 pixels.

Benefit of using mask

	Proposed	Pix2Pix	CycleGAN	StarGAN	GANimation
Size ($\times 10^4$)	1.68	4100	780	850	850
Time (msec)	2.70	320	710	580	507

- Comparison of MR with 4 state-of-the-art GAN architectures
- MR has more than two orders of magnitude fewer number of parameters than each of these GANs.
- MR is more than two orders of magnitude faster in synthesizing an expression.

Masked Regression – error formulation

- The error function for Masked Regression can be written as

$$E^{\text{MR}}(W) = \frac{1}{2} \|(W \circ M)X^T - T^T\|_F^2 + \frac{\lambda_M}{2} \|W \circ M\|_F^2 \quad (5)$$

- Only those weights are learned for which $m_{ij} = 1$. The rest are fixed to 0.
- Closed-form solution cannot be obtained due to the Hadamard product.

Masked Regression – error formulation

- Per-pixel decomposition

$$E^{\text{MR}}(W) = \sum_{i=1}^K E^{\text{MR}}(W^i) \quad (6)$$

where

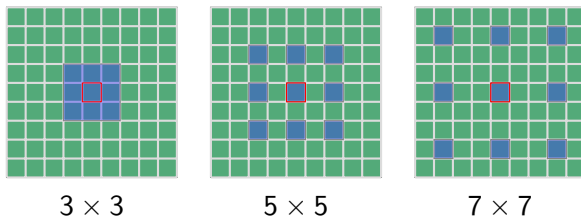
$$E^{\text{MR}}(W^i) = \frac{1}{2} \|(W^i \circ M^i)X^T - T_i^T\|_2^2 + \frac{\lambda_M}{2} \|W^i \circ M^i\|_2^2 \quad (7)$$

where W^i is the i -th row of W .

- Gradient and Hessian computations are a bit involved (refer to paper (Khan et al. 2019)).
- **Globally optimal W^i can now be computed in closed-form.**

Dilated receptive fields

- The proposed method can be easily modified to have not-so-local receptive fields.
- We use dilated receptive fields to observe larger input regions using the same number of weights.
- This helps to avoid over-fitting by limiting the complexity of the model.



Local vs Sparse

- Local receptive fields can be viewed as extremely sparse receptive fields with manually designed and fixed localizations.
- Alternative: learn sparse receptive fields.
- **Will a sparsely learned topology also converge to our local receptive fields?**

Local vs Sparse

- To answer this question we learn the receptive field W^i for each output pixel by minimizing the ℓ_1 -regularized sum of squared errors

$$\min_{W^i} \frac{1}{2} \|XW^i - T_i\|_2^2 + \lambda_1 \|W^i\|_1 \quad (8)$$

using the LASSO algorithm².

- We also learn by minimizing the ℓ_0 -regularized sum of squared errors

$$\min_{W^i} \frac{1}{2} \|XW^i - T_i\|_2^2 \text{ s.t. } \|W^i\|_0 \leq \lambda_0 \quad (9)$$

using the OMP algorithm³.

²Tibshirani 1996.

³Pati, Rezaiifar, and Krishnaprasad 1993.

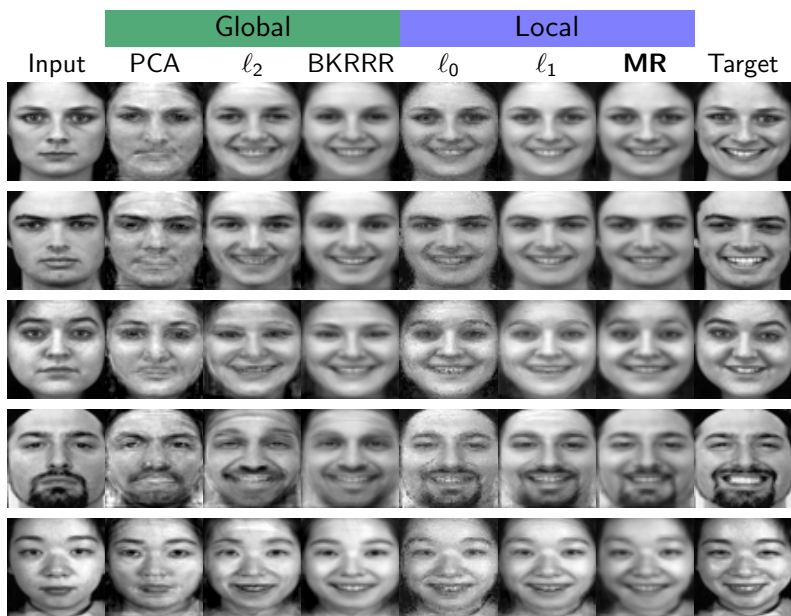


Figure: MR successfully synthesized a happy expression while preserving identity and retaining facial details the most.



Figure: For each neutral input, MR effectively transformed into 6 different expressions while preserving identities and facial details.

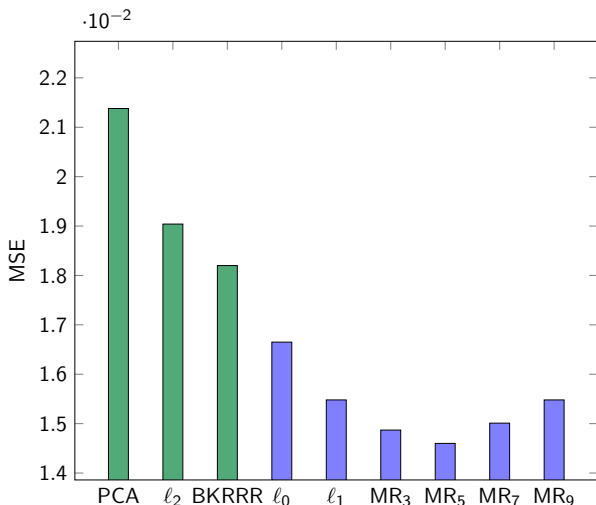


Figure: MSE for different methods averaged over 12 expression mappings. Employing too large a receptive field increased the MSE since long-range receptive fields fail to capture the local nature of facial expressions.

Learned Biases

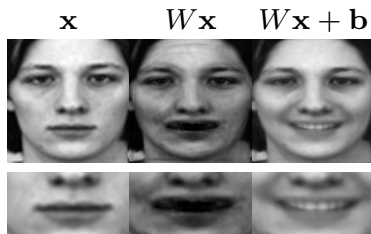
Grayscale



Color



Role of weights and biases



Grayscale synthesis

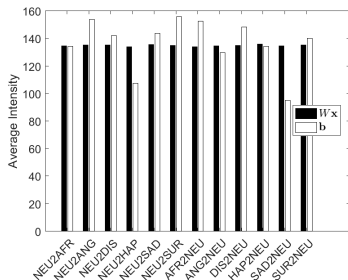


Color synthesis

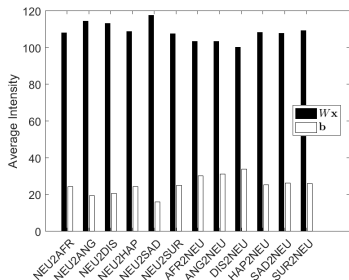
- Weights are predominantly used to transform the visible parts of the input expression into the target.
- Biases are used to insert hidden information such as teeth for a happy expression.

Over 12 expression mappings, we compare the average absolute intensity of the transformation produced by the weights with the additive transformation learned as biases.

ℓ_2 -regression



Masked regression



Bias often dominated the weights. Weights ~ 5 times as important as bias.
 Leads to loss of identity. Leads to better identity preservation.

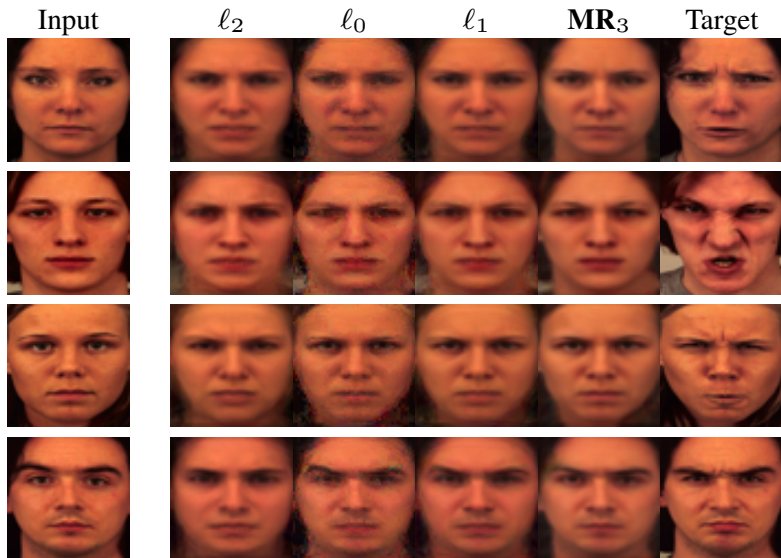


Figure: MR preserves background and other details unrelated to the desired expression.

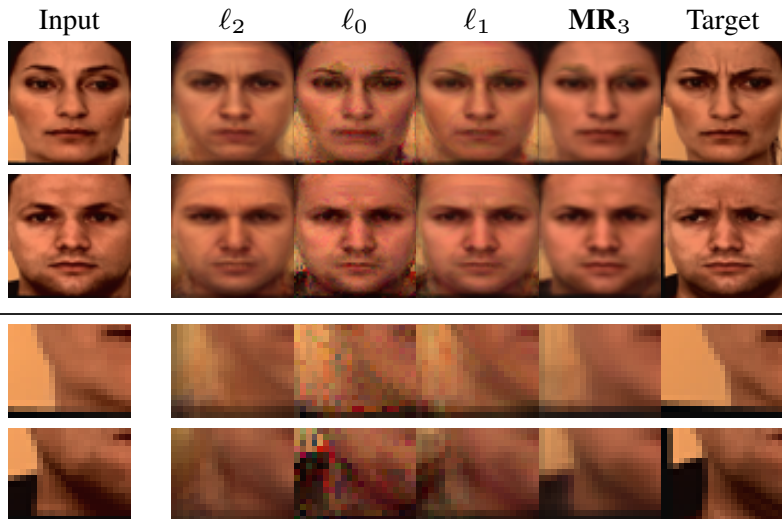


Figure: MR preserves background and other details unrelated to the desired expression.

MR for color images

$$E^{\text{CMR}}(W) = \frac{1}{2} \sum_{c=1}^C \|(W \circ M)X_c^T - T_c^T\|_F^2 + \frac{\lambda_M}{2} \|W \circ M\|_F^2 \quad (10)$$

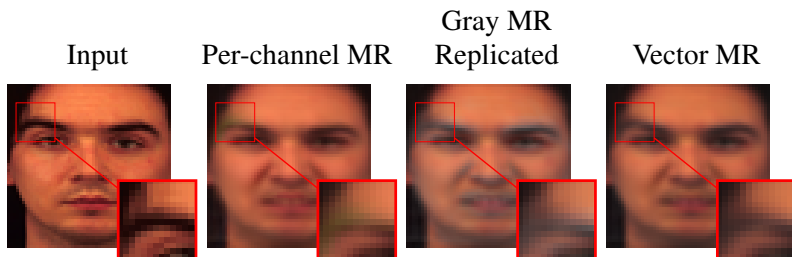


Figure: Regression on color tuples leads to lesser color leakage compared to separate regressions on each channel.

MR on non-frontal faces

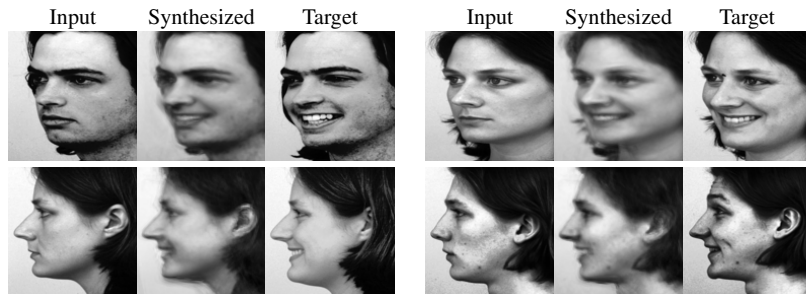


Figure: FES on non-frontal faces.

Out-of-dataset generalization

- MR is most effective for *out-of dataset* images.
- Images not belonging to any of the datasets used for training, validation and testing.
- Such images belong to significantly different distributions compared to training distribution.
- Three categories
 - ① People
 - ② Sketch drawings
 - ③ Animals



MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



Disgusted



Happy



Sad



Surprised



MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



Disgusted



Happy



Sad



Surprised





MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



Disgusted



Happy



Sad



Surprised



MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



Disgusted



Happy



Sad



Surprised





MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



Disgusted



Happy



Sad



Surprised



MR

 ℓ_1 ℓ_0 ℓ_2

Afraid



Angry



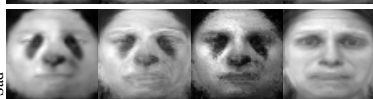
Disgusted



Happy



Sad



Surprised



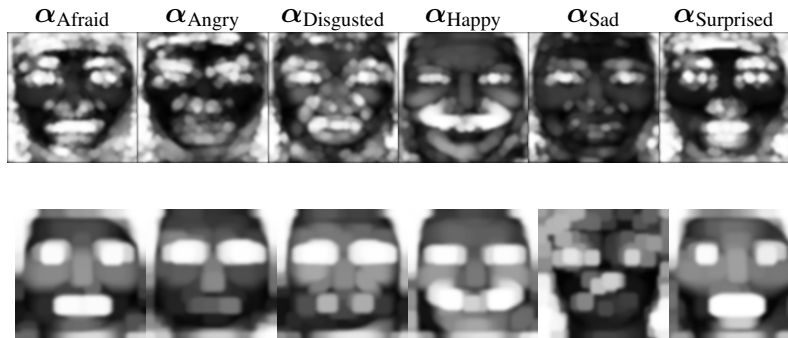


Figure: Top: Role of each pixel for expression generation. Higher intensity implies greater role. The role of the i -th pixel is computed entirely from its learned receptive field W^i . **Bottom:** Using different dilation and post-processing parameters.

$$\mathbf{y}' = (1 - \alpha) \circ \mathbf{x} + \alpha \circ \mathbf{y} \quad (11)$$

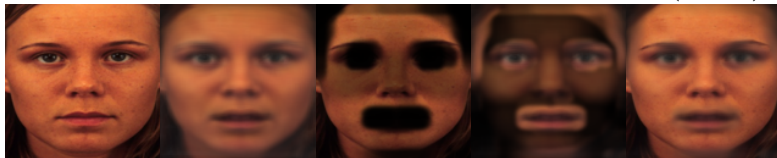
Input
 x

MR
 y

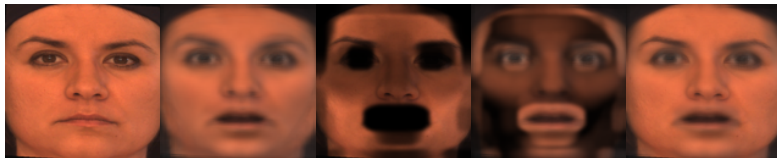
$(1 - \alpha) \circ x$

$\alpha \circ y$

RMR
 $(1 - \alpha) \circ x$
 $+ (\alpha \circ y)$



Afraid



Surprised



Angry

Input
 \mathbf{x}

MR
 \mathbf{y}

$(1 - \alpha) \circ \mathbf{x}$

$\alpha \circ \mathbf{y}$

RMR
 $(1 - \alpha) \circ \mathbf{x}$
 $+ (\alpha \circ \mathbf{y})$



Angry



Disgusted

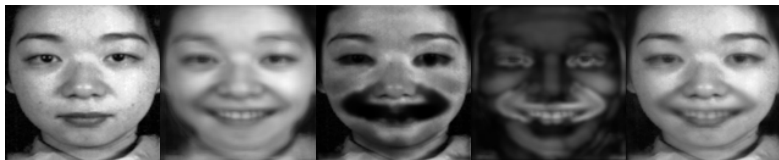
Input
 \mathbf{x}

MR
 \mathbf{y}

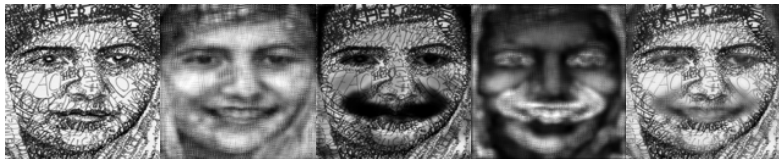
$(1 - \alpha) \circ \mathbf{x}$

$\alpha \circ \mathbf{y}$

RMR
 $(1 - \alpha) \circ \mathbf{x}$
 $+(\alpha \circ \mathbf{y})$



Happy



Happy

In-dataset images

Input



MR



RMR



Pix2Pix



CycleGAN



StarGAN



GANimation



Out-of-dataset images (sketches and animals)



Pix2Pix



CycleGAN



StarGAN



GANimation

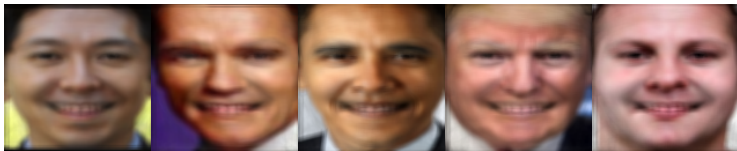


Out-of-dataset images (people)

Input



MR



Pix2Pix



CycleGAN



StarGAN



GANimation



Input

Afraid

Angry

Disgusted

Happy

Sad

Surprised



RMR



GANimation



RMR



GANimation



RMR

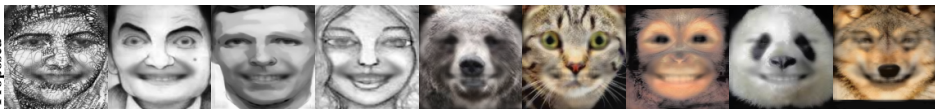


GANimation

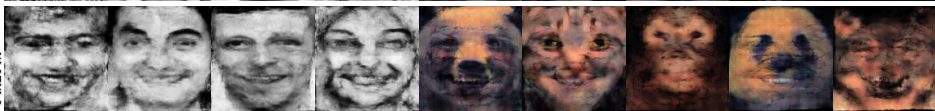
Input



Proposed



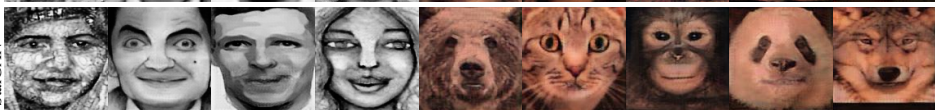
Pix2Pix



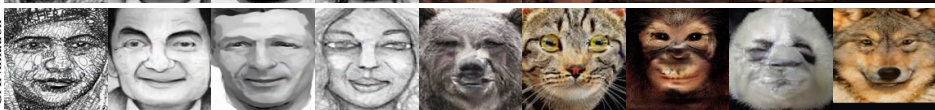
CycleGAN



StarGAN



GANimation



Quantitative comparison with GANs

Table: Drop in expression recognition accuracy (in percentage points) when changing from test set images to out-of-dataset images.

Pix2Pix	CycleGAN	StarGAN	GANimation	MR
35.72	16.39	21.43	20.74	12.39

Conclusion

- Constrained version of ridge regression for local receptive fields.
- Efficient closed-form solution of global minimum.
- Excellent learning ability on very small datasets despite simplicity.
- Easy implementation and extremely fast training.
- Better generalization despite using small training datasets.
- Extremely small model size.
- Intuitive interpretation of receptive fields exploited to refine results.
- Better out-of-dataset generalization compared to state-of-the-art GANs.

References II



N. Khan et al. "Masked Linear Regression for Learning Local Receptive Fields for Facial Expression Synthesis". In: *International Journal of Computer Vision (IJCV)* (2019).

References IV



Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad.
“Orthogonal matching pursuit: recursive function
approximation with applications to wavelet
decomposition”. In: *Proceedings of 27th Asilomar
Conference on Signals, Systems and Computers*. Nov.
1993, 40–44 vol.1.



A. Pumarola et al. "GANimation: One-Shot Anatomically Consistent Facial Animation". In: *International Journal of Computer Vision (IJCV)* (2019).



Arman Savran et al. "Bosphorus database for 3D face analysis". In: *European Workshop on Biometrics and Identity Management*. Springer. 2008, pp. 47–56.

Questions?

