

Problem

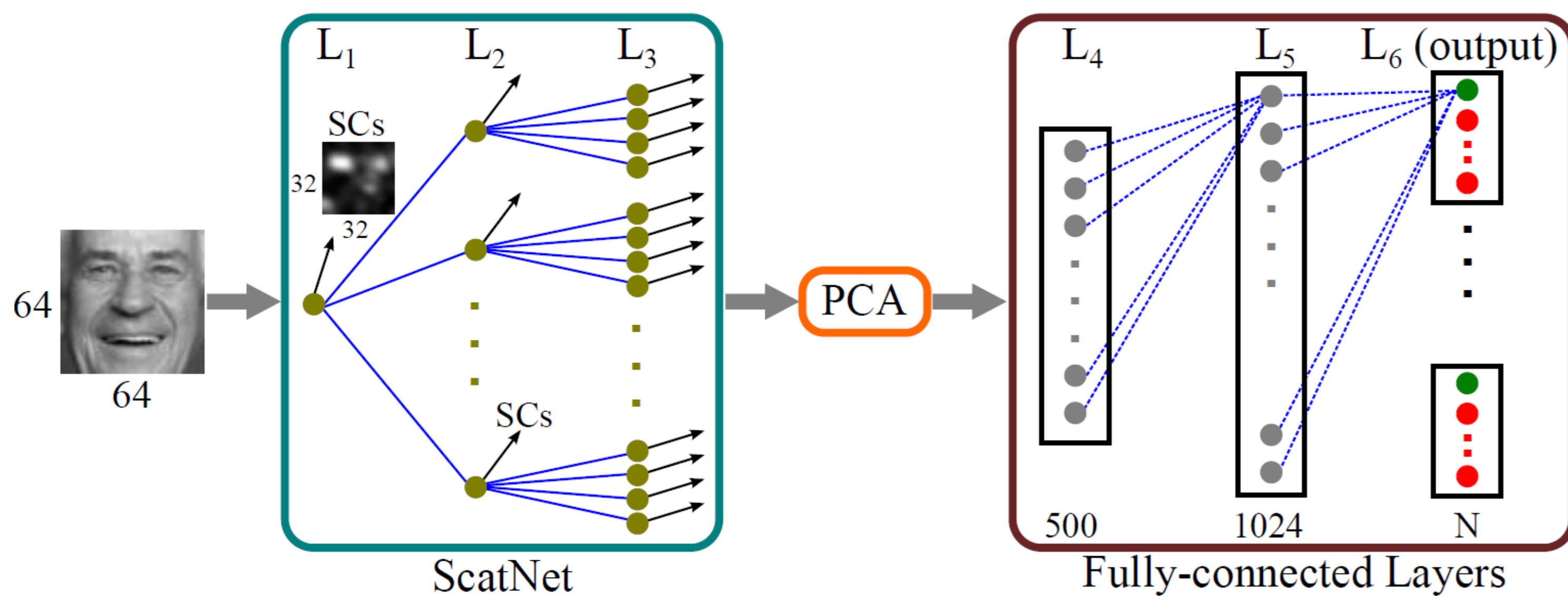
- Automatic age estimation (AAE) from face images is a challenging problem because of large facial appearance variations resulting from a number of factors, e.g., gender and facial expressions.
- We propose a generic, deep ranking model which first extracts features from the face through a scattering network (ScatNet), then reduces the feature dimension by principal component analysis, and finally predicts the age via category-wise rankers.

Contribution

- Develop a generic, deep ranking model that combines ScatNet and the ordinal ranking for AAE from face images.
- Our deep ranking model is general and can be applied to age estimation from faces with large facial appearance variations as a result of aging or facial expression changes.
- We show that the high-level concepts learned from large-scale neutral faces can be transferred to estimating ages from faces under expression changes, leading to improved performance.

Deep Ranker

- Our network model for human age prediction from face images



ScatNet as Facial Descriptors

- ScatNet [1] is a deep convolutional network of specific characteristics.
 - use predefined wavelets
 - compute by a cascade of wavelet transforms and modulus pooling operators from shallow to deep layers
- ScatNet can eliminate variability from translations, rotations or scaling.
- Unlike standard CNNs, ScatNet requires no learning of the parameters.

DeepRank: An Ordinal Regression Ranker

- Rank Encoding:** We employ the reduction framework [2] to conduct a deep ranker. Given a set of training samples $\mathbf{X} = \{(\mathbf{x}_i, \mathbf{y}_i), i = 1, \dots, N\}$, let $\mathbf{x}_i \in \mathbf{R}^D$ be the input image and \mathbf{y}_i be a rank label ($\mathbf{y}_i \in \{1, \dots, K\}$), respectively, where K is the number of age ranks. For rank k , we separate \mathbf{X} into two subsets, \mathbf{X}_k^+ and \mathbf{X}_k^- , as follows:

$$\begin{aligned} \mathbf{X}_k^+ &= \{(\mathbf{x}_i, +1) | \mathbf{y}_i > k\} \\ \mathbf{X}_k^- &= \{(\mathbf{x}_i, -1) | \mathbf{y}_i \leq k\}. \end{aligned} \quad (1)$$

- $K - 1$ such binary classifiers are required for K age ranks.
- The output layer can be encoded in an intuitive-to-understand structure: For a face image with true age k , the teaching vector with length $K - 1$ is designed as $[1, \dots, 1, -1, \dots, -1]$, where the first $k - 1$ values are 1 and the remaining -1 .
- Prediction:** Aggregation the network's binary outputs into a rank $r(\mathbf{x}_i)$:

$$\mathbf{r}(\mathbf{x}_i) = \mathbf{1} + \sum_{k=1}^{K-1} \mathbb{I}[\mathbf{O}_k(\mathbf{x}_i) > \mathbf{0}], \quad (2)$$

where $\mathbf{O}_k(\mathbf{x}_i)$ is the output of the k th node, and $\mathbb{I}[\cdot]$ is 1 if the condition is met and 0 otherwise.

DeepRank+: A Multi-task Ranker

- Category-wise Label Encoding:** Assume there are C category-wise rankers. The encoding for each ranker consists of two constituents: the category label(s) and the age rank.
- For category j , its encoding is given as $\mathbf{e}_j = [\mathbf{g}_j, \mathbf{r}_j]$, where \mathbf{g}_j denotes the desired output of the category and \mathbf{r}_j . Concatenating the C encoding sets forms a final encoding: $\mathbf{E} = [\mathbf{g}_0, \mathbf{r}_0, \dots, \mathbf{g}_j, \mathbf{r}_j, \dots, \mathbf{g}_C, \mathbf{r}_C]$.
- When a face image with age k belongs to category j , the teaching vector is defined as:

$$\left[\begin{array}{c} \mathbf{g}_1 \\ -1, \underbrace{0, \dots, 0}_{K-1}, \dots, \mathbf{g}_j \\ 1, \dots, 1, \underbrace{-1, \dots, -1}_{K-k}, \dots, \mathbf{g}_C \\ -1, \underbrace{0, \dots, 0}_{K-1} \end{array} \right], \quad (3)$$

where $\mathbf{0}$ s denote "don't care".

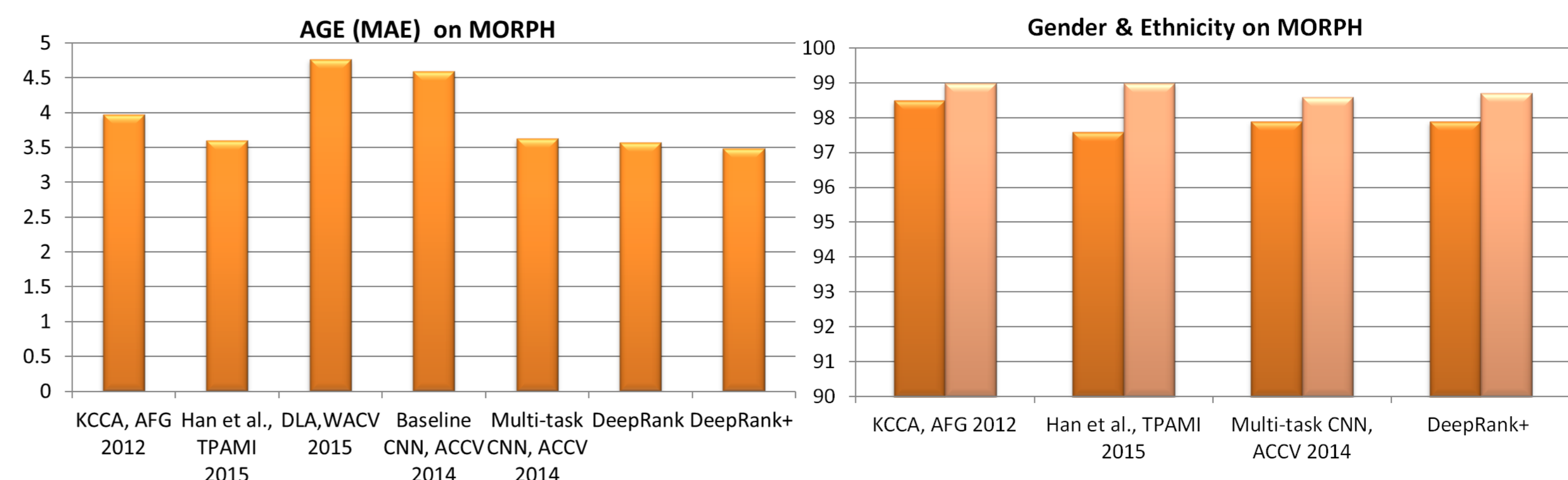
- Prediction:** First determine to which category the face image belongs. Then, we use Equation (2) to aggregate the ranking results of the determined category.

Experimental Results

- We conducted two sets of experiments on three face datasets.
 - Age estimation with different gender and races on a large-scale MORPH dataset.
 - Age estimation under expression changes on the Lifespan and FACES datasets.

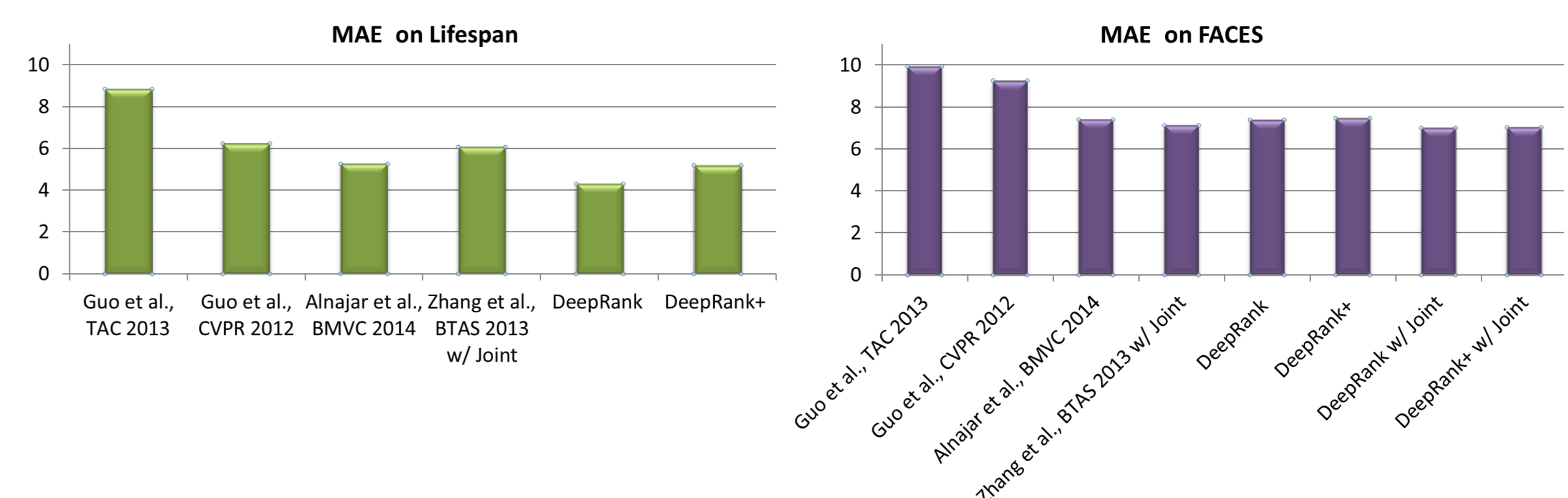
- Age Estimation with Different Gender and Races on MORPH

- MORPH** comprises more than 55,000 face images of 13,000 individuals aged from 16 to 77 years (and thus there are 62 age ranks).
- Joint learning of age, gender, and race not only boosts the performance on age estimation but also allows our network to accurately obtain other facial attributes from a face image.
- Network Models on MORPH: combinations of ethnicity (Black (B) and White (W)) and gender (Female (F) and Male (M)) the category labels, which results in 4 categories. Layer L_6 consists of 248 (i.e., $(1 + 61) \times 4$) nodes.



- Age Estimation under Expression Changes on Lifespan & FACES

- Lifespan** contains 1,046 images of 580 individuals in various expressions (e.g., anger, annoyed, disgusted, grumpy, happy, neutral, sad, and surprise), and the ages range from 18 to 93.
- Network Models on Lifespan: the number of nodes in L_6 of DeepRank is 75, and the number of nodes in L_6 of DeepRank+ is set to 152 (i.e., $(1 + 75) \times 2$).
- FACES** contains 2,052 images of 171 subjects ranging from ages 19 to 80 in two sets of six facial expressions (neutrality, sadness, happiness, disgust, fear, and anger).
- Network Models on FACES: DeepRank has 61 nodes in L_6 (62 age groups), and DeepRank+ has 372 nodes because 62 (1 category and 61 rank labels) nodes are required for each category-wise ranker and there are 6 expressions.
- Network Models on FACES with Extended Data: DeepRank has 75 nodes in L_6 (76 age groups) with additional data from Lifespan, and DeepRank+ has 456 (i.e., $(1$ (category) + 75 (ranks)) \times 6 (expressions)) nodes.



[1] Joan Bruna and Stéphane Mallat, Invariant scattering convolution networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2013.

[2] H.-T. Lin and L. Li, Reduction from cost-sensitive ordinal ranking to weighted binary classification. *Neural Computation*, 2012.