

SUPPLEMENTAL MATERIAL FOR LEARNING LABEL ENCODINGS FOR DEEP REGRESSION

Deval Shah & Tor M. Aamodt

Department of Electrical and Computer Engineering
University of British Columbia, Vancouver, BC, Canada
{devalshah, aamodt}@ece.ubc.ca

A APPENDIX

This supplemental material provides additional results and ablation studies (Section A.1), methodology for baseline encodings design approaches (Section A.2), and related work on task-specialized approaches and experimental setup (Section A.3) for RLEL. Code is available at https://github.com/ubc-aamodt-group/RLEL_regression

A.1 ABLATION STUDY

Section A.1.1, Section A.1.2, and Section A.1.3 provide an ablation study and supporting data on impact of proposed regularization functions and hyperparameters on label encoding learning. Section A.1.4 covers an ablation study on the impact of the number of fully connected layers in direct regression and multiclass classification. Section A.1.5 explains and compares deep hashing approaches (adapted for regression) with RLEL. Section A.1.6 provides results for geometric mean and Pearson coefficient as evaluation metrics.

A.1.1 IMPACT OF REGULARIZER R1

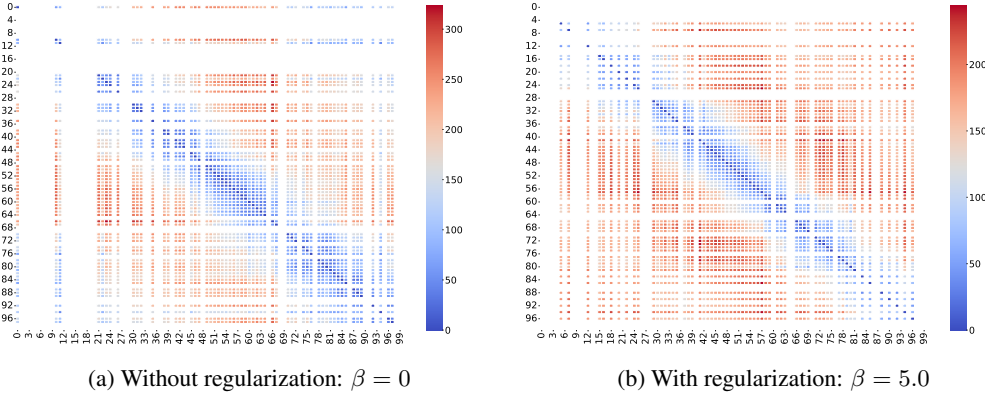


Figure 4: (a) and (b) show the L1 distance between pairs of encodings for FLD1_s benchmark for $\beta = 0$ and $\beta = 5.0$, respectively. Each cell (i, j) in this matrix represents the L1 distance between learned encodings for label i and j , i.e., $\|E_{i,:} - E_{j,:}\|_1$.

We proposed regularization function R1 to encourage the L1 distance between encodings to be proportional to the difference between corresponding label values. Figure 4a and Figure 4b represent the L1 distance between pairs of learned encodings for FLD1_s benchmark without and with regularization, respectively. The X-axis and Y-axis represent the label values. Here, some columns and rows are replaced by white lines, as these label values are not present in the training dataset. The data point at coordinates (i, j) represent the L1 distance between encodings for label i and j , i.e., $\|E_{i,:} - E_{j,:}\|_1$. For example, in Figure 4a the L1 distance between encodings for label values 0 and 97 is ~ 120 (light-blue coloured point at coordinate (0, 97)). In Figure 4b, the L1 distance between encodings for label values 4 and 96 is ~ 170 (red coloured point at coordinate (4, 96)).

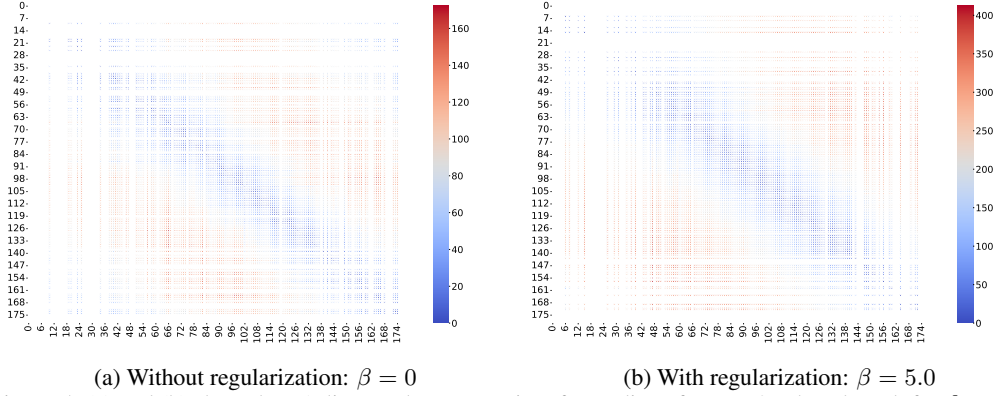


Figure 5: (a) and (b) show the L1 distance between pairs of encodings for FLD2_s benchmark for $\beta = 0$ and $\beta = 5.0$, respectively. Each cell (i,j) in this matrix represents the L1 distance between learned encodings for label i and j , i.e., $\|E_{i,:} - E_{j,:}\|_1$.

The first design property (Section 3) states that the L1 distance between encodings should increase with the difference between corresponding label values. The difference between label values for pairs of encodings increases with the distance from the diagonal of this plot. Thus, the value of data points (i.e., the L1 distance between encodings) should increase with the distance from the diagonal of this plot. As shown in Figure 4a without regularization, the distance between encodings is less for faraway label values (blue-colored data points away from diagonal), which shows that learned encodings do not follow the proposed design property. As shown in Figure 4b the introduction of regularization function R2 remedies this and increases the L1 distance between encodings for faraway labels. Similar observations are made for FLD2_s benchmarks, as shown in Figure 5a and Figure 5b

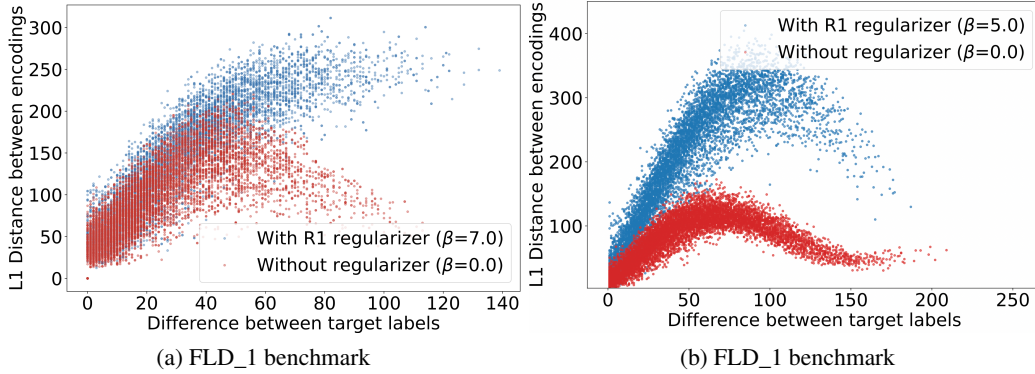


Figure 6: (a) and (b) plot the L1 distance between pairs of encodings versus distance between corresponding label values for FLD1_s and FLD2_s benchmarks.

Figure 6 plots the L1 distance between encodings versus the difference between corresponding label values for benchmarks FLD1_s and FLD2_s. For both the benchmarks, the proposed regularizer R1 helps enforce the first design property for real-valued label encodings and results in better label encodings with lower error (Table 3).

Effect of the scaling parameter in equation 3 We use the scaling parameter 2 in equation 3. Our intuition behind using the scaling parameter 2 is based on binary-encoded labels. For two adjacent labels (i.e., $|y_i - y_j| = 1$), the loss function encourages $\|\hat{Z}_i - \hat{Z}_j\|_1$ to be greater than 2. Here, \hat{Z} is the output encodings. In the case of binarized label encoding (-1 if $Z < 0$ and $+1$ if $Z > 0$), $\|Z_i - Z_j\|_1 = 2$ signifies that two encodings differ in at least one bit.

We also analyzed the effect of changing this parameter for two benchmarks. Table 6 shows the impact of changing this scaling parameter for two benchmarks. We observe that the error is higher if the scaling parameter is too low, as encodings for two adjacent labels can not be discriminated against. If this parameter is set too high, the encoding space is more constrained and consequently the performance is degraded.

Table 6: Effect of the scaling parameter on error for FLD1_s and FLD2_s benchmarks.

Value of the scaling parameter	NME (FLD1 _s)	NME (FLD2 _s)
1	4.89	4.15
2	4.71	4.15
3	4.83	4.20
4	4.97	4.27
5	4.95	4.28
6	5.06	4.41

Based on this intuition and empirical verification on two benchmarks, we use the value 2 for all benchmarks.

A.1.2 IMPACT OF REGULARIZER R2

The regularization function R2 is proposed to reduce the number of bit transitions in the learned label encoding. Figure 7 compares the label encodings learned for LFH1 benchmark for different values of α , where α is the weight of regularization function R2 (Equation 5). Each row k is the encoding for label value k . Each column k represents the output of the encoding position k for different label values. The regularization function is proposed to decrease the transitions in an encoding bit (blue \rightarrow red and red \rightarrow blue) over the range of label values. Section 4.4 provided quantitative results to demonstrate that increasing the value of α reduces the number of bit transitions. We observe similar trends in the plots of learned label encodings shown in Figure 7: increasing the value of α decreases bit transitions in the learned label encodings and improves MAE.

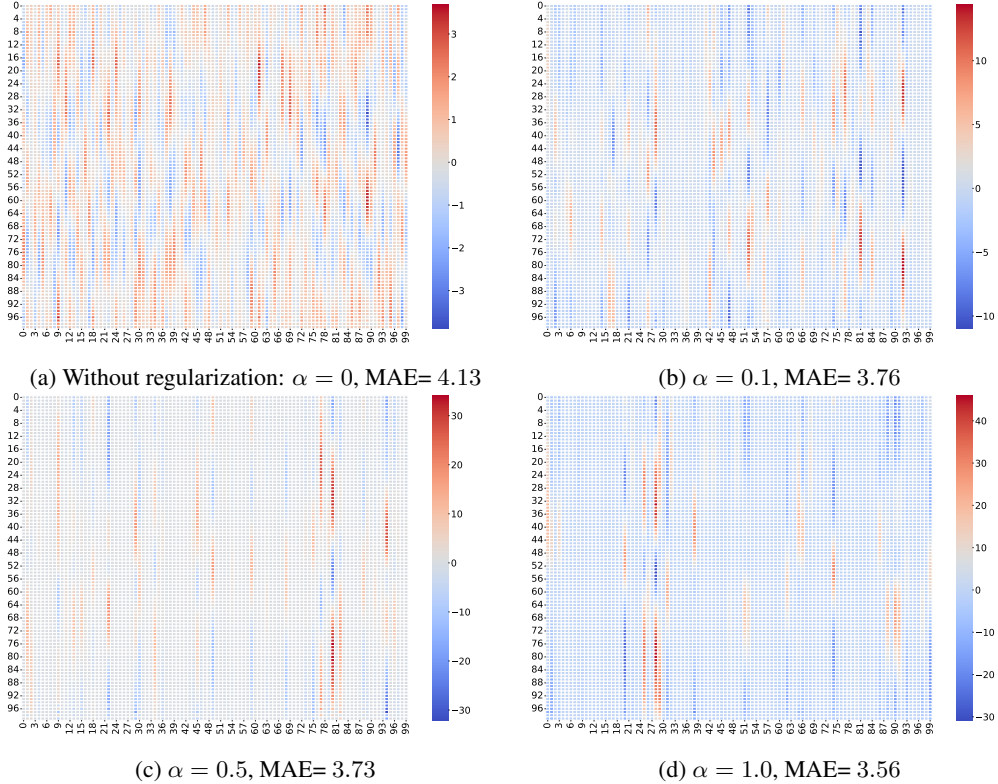


Figure 7: (a)-(d) represent the label encodings learned by RLEL for different values of weight α for regularizer R2 (Equation 5).

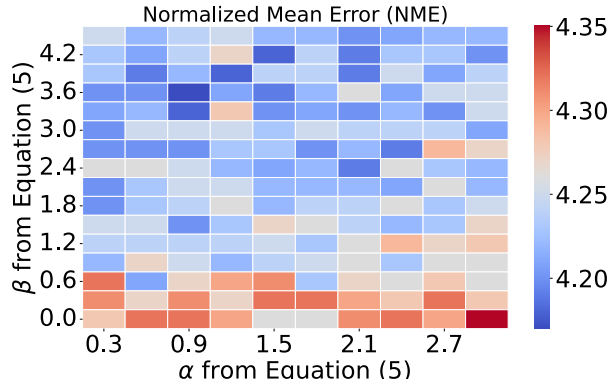


Figure 8: Impact of hyperparameters α and β from Equation 5 on NME for FLD1_s benchmark.

Table 7: Impact of the number of fully-connected layers in direct regression and multiclass classification on the error (MAE or NME). This table is reproduced from (Shah et al., 2022).

Benchmark	Direct regression		Multiclass classification	
	1 FC layer	2 FC layers	1 FC layer	2 FC layers
LFH1	4.76	5.19	4.49	4.82
LFH2	5.65	5.59	5.31	5.42
FLD1	3.60	3.63	3.58	3.56
FLD2	3.54	3.58	3.51	3.62
FLD3	4.64	4.63	4.50	4.64
FLD4	1.51	1.51	1.56	1.53
AE1	2.44	2.35	2.75	2.81
AE2	3.21	3.14	3.38	3.40
PN	4.24	4.33	4.56	5.74

A.1.3 EFFECT OF HYPERPARAMETERS IN RLEL

The RLEL approach introduces two hyperparameters. We first evaluate the sensitivity to these hyperparameters to determine the complexity of hyperparameter tuning. Figure 8 shows the NME for FLD1_s benchmark for different values of α and β values in Equation 5. As shown in the figure, the error is not sensitive to small changes in these hyperparameters' values, suggesting that a sparse search in the hyperparameter space suffices. Furthermore, several approaches have been proposed for efficient hyperparameter search (Li et al., 2017; Falkner et al., 2018), and any off-the-shelf hyperparameter tuners/libraries can be used to automatically find these values without manual efforts. In contrast, hand-designed codes need human intervention to design codes. Also, multiple training runs are still required to find suitable codes for a given benchmark from a set of hand-designed codes. On the other hand, RLEL provides an end-to-end automated approach for label encoding learning.

A.1.4 IMPACT OF THE NUMBER OF FULLY-CONNECTED LAYERS:

For RLEL, we use an extra fully connected bottleneck layer in the regressor as proposed by the prior work on regression by binary classification (Shah et al., 2022). We provide an ablation study (reproduced from (Shah et al., 2022)) to show the impact of additional fully connected layers in direct regression and multiclass classification. Table 7 provides the error (MAE or NME) for direct regression and multiclass classification with one or two fully connected layers after the feature extractor. As shown in the table, increasing the number of fully connected layers in direct regression and multiclass classification does not reduce the error for most benchmarks (possibly due to overparameterization).

A.1.5 COMPARISON WITH DEEP HASHING APPROACHES

Deep supervised hashing approaches use neural networks as a hash function and learn hash codes in an end-to-end manner. The loss function for deep supervised hashing is designed to preserve the similarity between inputs in the hashing space. Often, these approaches use the label information to determine the similarity between images (i.e., same label) (Liu et al., 2016; Xia et al., 2014). Some deep hashing approaches have proposed to augment the loss function with classification loss

to improve the performance. We adapt two widely used deep-hashing approaches to regression and compare RLEL with deep hashing approaches.

Liu et al. (2016) proposed a deep supervised hashing (DSH) approach with a loss function based on the pairwise similarity between images. The proposed approach introduces a loss function to preserve the similarity between output codes for similar training images (e.g., same class) and maximize discriminability between output codes for different training images (e.g., different class). Further, they propose using relaxation on the binary output and a regularizer to encourage the output code to be close to discrete values $+1/-1$. The hamming distance between output codes can be computed for binary-like outputs using the L2 norm. We use DSH for regression with some modifications (DSH-reg). We used the quantized label to determine the class of a training sample.

Lai et al. (2015) proposed a triplet ranking loss to learn a hash function that preserves relative similarities between images (SFLH). For images $(I, I+, I-)$, where I is closer to $I+$ than $I-$, the loss function is designed to encourage higher hamming distance between codes for $(I, I-)$ than $(I, I+)$. For classification datasets, triplets are typically formed using two images from the same class and one from a different class (Norouzi et al., 2012). They proposed to use a piece-wise threshold function to encourage binary-like outputs.

We use the above approach (SFLH) for regression with a few modifications (SFLH-reg). To generate triplets, we pick sets of three images from a given batch and determine the similarity between images using differences between the label values. We use K^2 triplets for a minibatch of K training samples.

Further, for both DSH-reg and SFLH-reg, we augment the loss function with regression loss. We add a fully-connected layer between the output code and prediction. The MSE loss between the final outputs and target labels is added to the loss function (DSH-reg-L2, SFLH-reg-L2).

Table 8: Comparison of RLEL with different deep hashing approaches adapted for regression.

Method	MAE
DSH-reg	71.3
DSH-reg-L2	4.11
SFLH-reg	69.8
SFLH-reg-L2	4.73
RLEL (only R1)	3.93
RLEL (R1 + R2)	3.55

Table 8 compares the modified deep hashing approaches with RLEL. The gap between loss functions with and without regression loss is significant, which shows that a loss function that only aims to preserve the similarity between output codes is not sufficient and needs to account for the error between decoded output and target (i.e., regression loss). RLEL results in a lower error as it is designed for regression problems that account for classifiers’ nonuniform error probability distribution.

Regularizer R1 encourages the distance between output codes for images to be proportional to the difference between label values, similar to pairwise or ranking-based loss functions proposed by deep hashing. However, deep hashing approaches use the hamming distance between binary outputs. As we show in Section 3.2, the hamming distance between codes does not account for the error probability of classifiers. Thus we use the L1 distance between the real-valued outputs to account for the confidence of the classifiers. R1 does not use regularizer or nonlinear activation on the output codes to encourage binary-like outputs, as typically done in deep hashing approaches. In contrast, we show that suitable regression codes can be learned by not using this constraint. Thus RLEL with only R1 regularizer results in lower error than deep hashing approaches.

A.1.6 EVALUATION

Table 9 compares RLEL with direct regression and multiclass classification using geometric mean and Pearson coefficient as evaluation metrics. The geometric mean represents the geometric mean of absolute error for the test dataset. The Pearson coefficient represents the correlation between the target and predicted labels for the test dataset. As shown in the table, RLEL results in significant reduction in the error compared to other generic regression approaches.

Table 9: Comparison of RLEL with different regression approaches using Geometric mean and Pearson coefficient as evaluation metrics.

	RLEL		Direct Regression		Multiclass Classification	
	GeoMean	Pearson Coeff.	GeoMean	Pearson Coeff.	GeoMean	Pearson Coeff.
LFH1	1.95	97.68	2.91	97.10	2.30	94.60
LFH2	2.09	92.22	2.49	91.06	2.40	88.76
FLD1	0.96	99.94	1.07	99.93	1.04	99.93
FLD1_s	1.31	99.87	6.38	99.81	1.81	99.80
FLD2	1.92	99.97	2.12	99.97	2.07	99.97
FLD2_s	2.44	99.96	3.03	99.98	3.22	99.94
FLD3	0.96	99.99	1.05	99.99	1.01	99.97
FLD3_s	1.21	99.98	1.56	99.97	1.37	99.97

A.1.7 THEORETICAL ANALYSIS OF PROPOSED REGULARIZATION FUNCTIONS

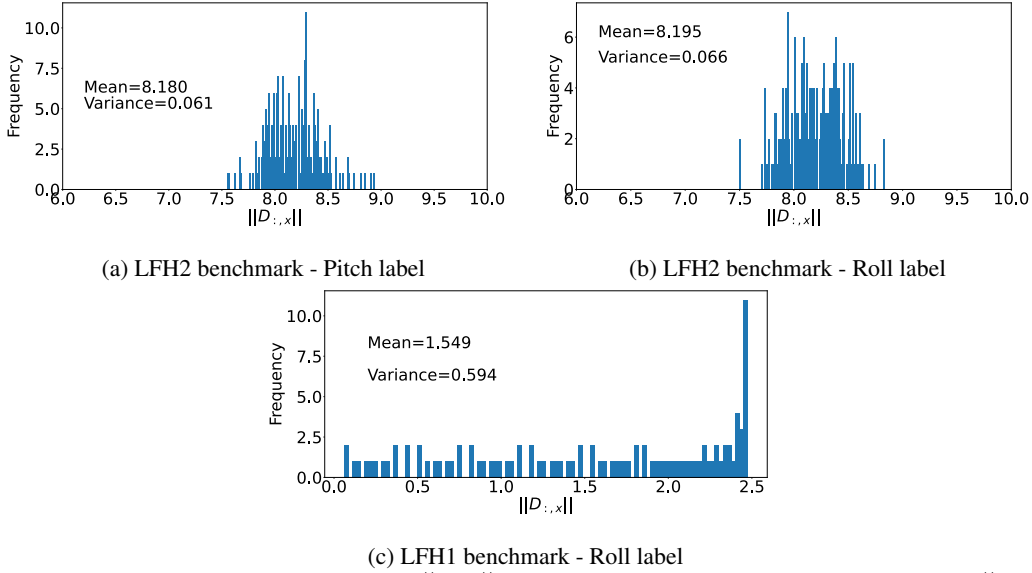


Figure 9: (a) and (b) plot the distribution of $\|D_{:,x}\|$ for LFH2 benchmark. (c) plots the distribution of $\|D_{:,x}\|$ for LFH1 benchmark. Here the variance is very low, which suggests that the assumption $\|D_{:,x}\| \approx \|D_{:,y}\|$, $x \in [1, N]$, $y \in [1, N]$ is valid. For the LFH1 benchmark, the variance is higher than LFH2. However, all outliers are for label values with very few (or even zero) training examples.

Regularization function R2:

We used matrix D instead of label encoding E to apply regularizer R2 in equation 4. We insight into this decision as follows. First, note the output encodings are multiplied with D to generate the correlation vector \hat{C}_i (Figure 3). We use the multiclass classification loss between \hat{C}_i and the target labels for training. Due to this, label encoding E and decoding matrix D are related, and use of matrix D proves to be effective for regularizer R2. We further explain this in detail below.

Let E represent an encoding matrix of size $N \times M$. Each row $E_{k,:}$ represents the encoding output when the label is k . D is the decoding matrix of size $M \times N$. Let \hat{C}_k represent the output correlation row vector of size $1 \times N$ when the target label is k . Here, \hat{C}_k is obtained by multiplying $E_{k,:}$ with D (Figure 3).

$$\hat{C}_k = E_{k,:} D \quad (6)$$

Since we apply softmax on the output vector to find the predicted label (Figure 3), ideally, \hat{C}_k^k should have the highest value as the target label value is k .

$\therefore \hat{C}_k^k > \hat{C}_k^x$, where, $x \neq k, x \in \{1, 2, \dots, N\}$

$\therefore E_{k,:} \cdot D_{:,k} > E_{k,:} \cdot D_{:,x}$, where, $x \neq k, x \in \{1, 2, \dots, N\}$ (Using equation 6). Let $\theta_{k,x}$ represent the angle between row vector $E_{k,:}$ and column vector $D_{:,x}$. This leads to the below equation:

$$||E_{k,:}|| ||D_{:,k}|| \cos(\theta_{k,k}) > ||E_{k,:}|| ||D_{:,x}|| \cos(\theta_{k,x}), \text{ where, } x \neq k, x \in \{1, 2, \dots, N\} \quad (7)$$

Shah et al. (2022) used a hand-crafted decoding matrix D with an equal number of 1s and 0s in each column for binary-encoded labels. Hence the L2 norm of each column is the same. In label encoding learning, parameters of matrix D are learned during training and are not constrained to have the same L2 norm for each column. However, we observe a similar trend empirically. Figure 9c plots the distribution of $||D_{:,x}||$ for different benchmarks. As shown in the figure, for most benchmarks, we observe a small variance in the distribution of $||D_{:,x}||$. Based on this intuition and empirical validation, we assume that $||D_{:,x}|| \approx ||D_{:,y}||$ for $x \in [1, N]$ and $y \in [1, N]$ to simplify the analysis.

Using this assumption in equation 7 leads to the following inequality:

$$\cos(\theta_{k,k}) > \cos(\theta_{k,x}), \text{ where, } x \neq k, x \in \{1, 2, \dots, N\}$$

Thus the cosine similarity between $E_{k,:}$ and $D_{:,k}$ should be the highest to predict the label k . The optimization process to reduce the loss between the target and prediction will try to maximize this cosine similarity. In the best case, the angle between $E_{k,:}$ and $D_{:,k}$ will be zero, and both vectors are parallel.

This simplification leads to the following relation between E and D .

$$E_{k,:} = t D_{:,k}, \text{ where } t > 0$$

$$\text{Similarly, } E_{k+1,:} = t' D_{:,k+1}, \text{ where } t' > 0$$

Since t and t' both are positive values, reducing $D_{i,k} - D_{i,k+1}$ also reduces $E_{k,i} - E_{k+1,i}$.

Regularizer rule R2 proposes to regularize the number of decision boundaries by regularizing $\sum_{i=1}^M \sum_{j=1}^{N-1} |E_{j,i} - E_{j+1,i}|$ as per equation 2. Based on the analysis above, regularizing $\sum_{i=1}^M \sum_{j=1}^{N-1} |D_{i,j} - D_{i,j+1}|$ helps with the above goal as $E_{j,i} - E_{j+1,i}$ reduces with $D_{i,j} - D_{i,j+1}$.

Regularization function R1:

The first property suggests $||E_{i,:} - E_{j,:}||_1 \propto |i - j|$.

So ideally, $||E_{i,:} - E_{j,:}||_1 = \lambda |i - j|$

Since $E_{x,:}$ is average of \hat{Z}_i for samples with label value x (equation 1), the above condition leads to:

$$||\hat{Z}_i - \hat{Z}_j||_1 = \lambda |y_i - y_j| \quad (8)$$

Based on this requirement, we add a regularization function $\max(0, \lambda |y_i - y_j| - ||\hat{Z}_i - \hat{Z}_j||_1)$, which penalizes the encodings if $||\hat{Z}_i - \hat{Z}_j||_1 < \lambda |y_i - y_j|$. It does not strictly impose equation 8. However, it approximately imposes the constraint as per shown in empirical verification in Section A.1.1

Our intuition behind using the scaling parameter 2 is based on binary-encoded labels. For two adjacent labels (i.e., $|y_i - y_j| = 1$), the loss function encourages $||\hat{Z}_i - \hat{Z}_j||_1$ to be greater than 2. Here, \hat{Z} is the output encodings. In the case of binarized label encoding (-1 if $Z < 0$ and $+1$ if $Z > 0$), $||Z_i - Z_j||_1 = 2$ signifies that two encodings differ in at least one bit.

Table 10: Impact of the number of quantization levels on error for FLD1 benchmark

Quantization levels (N)	NME
32	3.49
64	3.36
128	3.36
256	3.36
384	3.37
512	3.37

Table 11: Effect of dataset size on the error for FLD1 benchmark.

%Dataset used	RLEL	BEL	Difference (RLEL-BEL)
100	3.36	3.35	0.01
80	3.43	3.42	0.01
60	3.53	3.47	0.06
40	3.77	3.72	0.05
20	4.08	4.04	0.04
10	4.71	4.63	0.08

A.1.8 IMPACT OF THE NUMBER OF QUANTIZATION LEVELS (N)

The number of quantization buckets is treated as a design parameter for binary-encoded labels. [Shah et al. \(2022\)](#) showed that the error changes with the number of quantization levels. Fewer levels introduce quantization error, and more levels increase the number of bits in the encoding. They showed a trade-off between these two factors to decide the number of quantization levels.

Our work focuses on the design space of encoding and decoding functions. Hence we use the same values for the quantization levels (N) as BEL [Shah et al. \(2022\)](#). Parameter N tuning can be integrated into hyperparameter tuning or included in the optimization process.

We further analyze the effect of the number of quantization levels for RLEL. Table [10](#) shows the NME (Normalized Mean Error) for different values of N for FLD1 benchmark.

This suggests that the proposed method RLEL is less sensitive to the number of quantization levels for higher values. For RLEL, the decoding matrix that converts the encodings to the predicted label is also learned during the training (Figure 3). This matrix is of size $M \times N$, where each column represents the weight parameters for one quantization level. One possible reason for the above results is that matrix D learns the number of quantization levels suitable for this problem.

There is a potential to learn the number of quantization levels and non-uniform quantization using the proposed RLEL framework. For example, in Figure 3- step (4), fixed parameters j are used to scale the correlation vector \hat{C}_i^j and find the expected prediction \hat{y}_i . These parameters represent quantization levels. One possible approach to learning the quantization levels is to make these parameters trainable. In this case, L1/L2 loss between the expected prediction \hat{y}_i and target labels y_i can be used to train the network.

A.1.9 IMPACT OF DATASET SIZE ON ERROR FOR RLEL AND BEL

In order to compare the effect of dataset size on encoding design, we run BEL and RLEL approaches with the same training loss function (cross entropy loss in equation [5](#)). We take the dataset FLD1 and use a fraction of the dataset for training. The entire test dataset is used for testing here. Table [11](#) summarizes the error achieved by RLEL and BEL for different fractions of the training dataset. The evaluation shows that the gap between the performance of RLEL and BEL decreases with the increase in dataset size, which suggests that RLEL might be able to achieve lower error for larger datasets.

A.1.10 COMPARISON OF LEARNED AND MANUALLY DESIGNED ENCODINGS

We visually compare the encoding learned by RLEL with BEL manually designed code for one benchmark. Figure [10](#) shows the learned and manually designed encodings. Here, row k represents

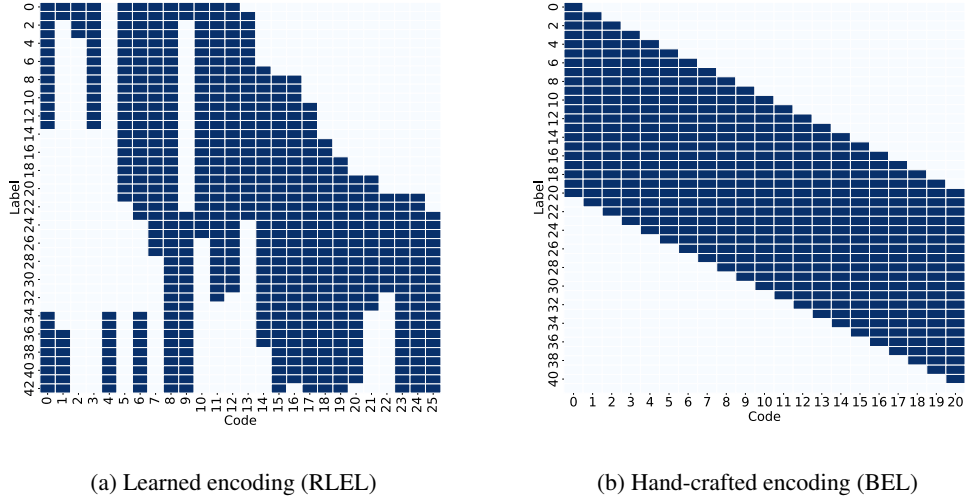


Figure 10: (a) and (b) give examples of learned and hand-crafted encodings. Here, row k represents an encoding for label k . Column j represents the bit values for classifier- k over the numeric range of labels.

an encoding for label k . Column j represents the bit values for classifier- k over the numeric range of labels. We notice some common characteristics between both encodings. For example, the codes for nearby labels differ by fewer bits than faraway labels. Both the codes also have fewer bit transitions ($0 \rightarrow 1$ and $1 \rightarrow 0$ transitions in a column). These characteristics in the learned encodings are encouraged by the proposed regularizers R1 and R2. There are a few differences between learned and hand-crafted encodings. In contrast to hand-crafted labels, encodings for adjacent labels do not differ in some cases, where hand-crafted encoding assures at least one or two bits of difference between adjacent labels.

A.2 LABEL ENCODING DESIGN

We evaluate different label encoding design approaches, including simulated annealing and autoencoder. These approaches have been used to design encodings for multiclass classification by prior works (Song et al., 2021; Cissé et al., 2012). We adapt these approaches to design encodings for regression tasks and compare RLEL with these code design techniques. This section provides the methodology for simulated annealing and autoencoder-based label encoding design.

A.2.1 SIMULATED ANNEALING

Simulated annealing is a probabilistic approach to find a global optimum of a given function. It is often used for combinatorial optimization, where the search space is discrete. Algorithm 1 represents the pseudo-code for label encoding design using simulated annealing. This algorithm takes two hyperparameters, K_{\max} (number of iterations) and T (initial temperature). It designs a code matrix C of size $N \times M$, where N is the number of values and M is the number of bits. Each row k in this code matrix represents encoding for value k . Code matrix C is initialized with a random matrix of 0 and 1 (Line 1).

For each iteration, a new code matrix C_{new} is sampled from the current code matrix C using a Move function (Line 4). For example, a move function can be designed to randomly flip a few bits in C . The difference between the errors of the current and new code matrix is measured (Line 5). The error of a code matrix, i.e., expected regression error for this problem, is measured using function E . For example, E can be replaced by training a regression network for a given code matrix to measure the regression error. Finally, the current code matrix C is updated with the new matrix C_{new} probabilistically. The probability is determined using the decrease in regression error and current temperature t (Line 6-8). The current temperature is updated for each iteration (Line 9).

Algorithm 1 Simulated annealing for encodings design**Input:** K_{\max}, T, M, N ;**Output:** $C \in \{0, 1\}^{M \times N}$;

```

1:  $C = C_0 \in \{0, 1\}^{N \times M}$ , where  $\Pr(C_{0,i,j} = 0) = \Pr(C_{0,i,j} = 1)$ 
2:  $t = T$ 
3: for  $k \in K_{\max}$  do
4:    $C_{\text{new}} = \text{Move}(C)$ 
5:    $D = E(C_{\text{new}}) - E(C)$ 
6:   if  $D < 0$  or  $e^{-\frac{D}{t}} > \text{Random}(0,1)$  then
7:      $C = C_{\text{new}}$ 
8:   end if
9:    $t = T / (k + 1)$ 
10: end for

```

There are mainly two design parameters in the above algorithm: the error measurement function E and the move function Move . We further explain the design of these functions.

Error measurement:

We used the expected absolute error between targets and decoded predictions for a given code matrix as its error, as the goal is to design a code matrix that results in lowest regression error. However, training a regression network for each sample code matrix to measure its regression error is computationally expensive and time-consuming (~ 200 training runs). Hence we use an analytical model to estimate the regression error for a given code matrix.

Regression error is the absolute error between targets Q_i and decoded predictions \hat{Q}_i . For a given target Q_i and target code $B_i = C_{Q_i,:}$, the predicted code (\hat{B}_i) will be erroneous due to classification errors. This erroneous predicted code is decoded to a predicted value (\hat{Q}_i). The following equation is used to predict \hat{Q}_i in expected-correlation-based decoding (Shah et al., 2022).

$$\mathcal{D}^{\text{GEN-EX}}(\hat{B}_i, C) = \sum_{k=1}^N k \sigma_k, \text{ where } \sigma_k = \frac{e^{\hat{B}_i \cdot C_{k,:}}}{\sum_{j=1}^N e^{\hat{B}_i \cdot C_{j,:}}} \quad (9)$$

The regression error can be estimated given sufficient samples of B_i and \hat{B}_i .

Shah et al. (2022) provided an approximate model of classification errors. They showed that for each classifier, its error probability distribution can be approximated using a combination of p Gaussian distributions, where p is the number of bit transitions. Each Gaussian distribution is centered around a bit transition. For example, for bit- k in unary code with bit transition between $Q = k$ and $Q = k + 1$, the error probability of the classifier- k for different target labels Q_i can be approximated as:

$$e_k(Q_i) = r f_{\mathcal{N}(\mu_k, \sigma^2)}(Q_i), \text{ where, } \mu_k = k + 0.5 \quad (10)$$

\hat{B}_i can be sampled for the given Q_i and C using the above error-probability model. Equation 9 is then used to find the decoded prediction \hat{Q}_i . We measure the expected absolute error between \hat{Q}_i and Q_i using $100 \times N$ samples.

We further verify the validity of this analytical model by finding the correlation between regression error measured by this model and trained networks. Figure 11 plots the analytical regression error versus actual regression error for FLD_1 benchmarks. Here, each point is for a different code matrix. The Y-axis represents the absolute error approximated by the proposed analytical model. The X-axis represents the absolute test error of a trained network for a given code matrix. The figure shows that the proposed analytical model for error measurement approximates error with significant speedup.

Move function: The move function flips some bits in the current code matrix to sample a new one. A naive approach would be to randomly flip b bits. We further optimize the move function to consider the regression task objective. For a given code matrix, using the proposed analytical model, we find

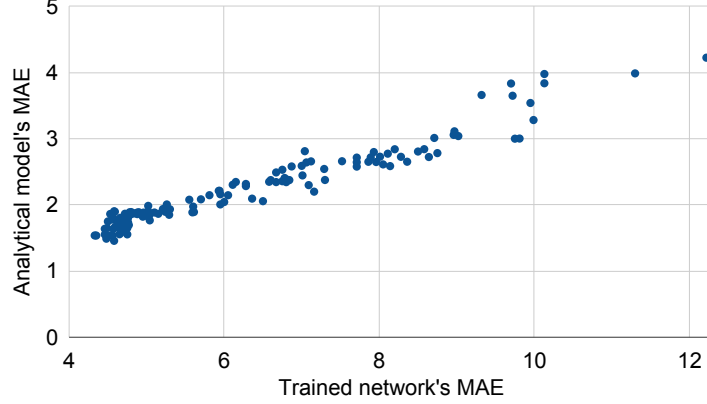


Figure 11: Comparison of Mean Absolute Error (MAE) approximated by the proposed analytical model and trained network for different code matrices. Each point in this plot is a different code matrix.

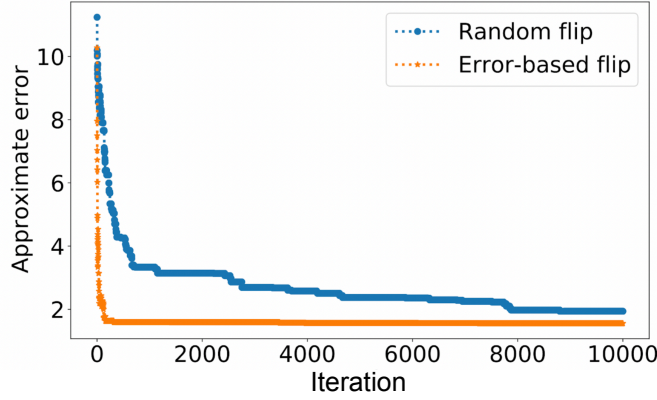


Figure 12: Comparison of convergence of random-flip and proposed error-based flip move functions.

a matrix F of size $N \times N$, where $F_{i,j} = |i - j| \times \Pr(\text{Round}(\mathcal{D}(\hat{B}, C)) = j | Q = i, B = C_{i,:})$. Thus, each element represents a pair $(C_{i,:}, C_{j,:})$ of encodings' contribution to expected error. We select top- b pairs from this matrix. For each pair of encodings, we find bit-positions that have equal bit-values between two encodings, and a randomly selected bit-position from this list is flipped in encoding $C_{i,:}$. This procedure increases the hamming distance between pairs of encodings that contribute the highest to the regression error. Figure 12 compares the convergence of the proposed move function and a random-flip-based move function. Here the Y-axis represents the approximated error for the current code matrix, and X-axis represents the iteration identifier. The figure shows that the proposed error-based move function results in faster convergence and lower error.

We use the proposed move function with the analytical model to approximate regression error in Algorithm 1 to design label encoding for regression using simulated annealing.

A.2.2 AUTOENCODER

Cissé et al. (Cissé et al., 2012) proposed an autoencoder-based approach to design encodings for a multiclass classification problem. Figure 13 represents the network architecture used for encodings design. Input S_i is an N -dimensional vector for class i . Here, each element $S_i[j]$ represents the similarity between class i and j . The output of the bottleneck layer C_i is the designed encodings for class i .

For regression problems, we set $S_i[j] = |i - j|$. Let W represent the weight parameters of the network. The network is trained using SGD optimization, where each batch consists of randomly

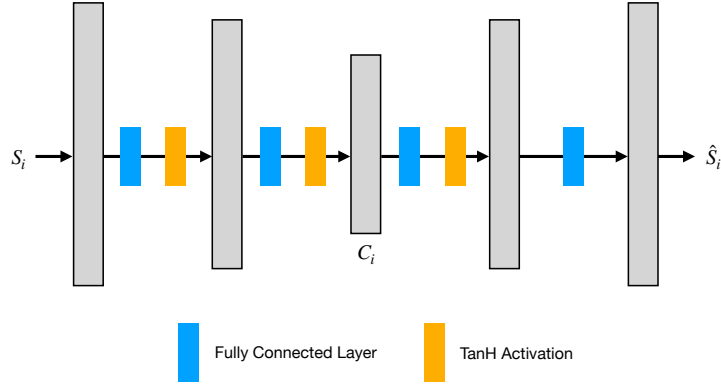


Figure 13: Network architecture for autoencoder-based encodings design.

sampled i and j . The following loss function is used for training:

$$\mathcal{L} = \|\hat{S}_i - S_i\|^2 + \|\hat{S}_j - S_j\|^2 + \beta \max(0, b - \|C_i - C_j\|_1) + \gamma \|W\|^2 \quad (11)$$

Here, the first and second terms represent reconstruction losses for inputs S_i and S_j . The third term encourages a minimum distance of b between any pair of encodings to yield unique encodings for different classes. The fourth term is an L2-regularizer.

Once the network is trained, the real-valued encodings C_i are converted to binary encodings such that it has equal numbers of 0s and 1s. This formulation introduces three hyperparameters. We determine the number of bit transitions in the designed label encodings and select hyperparameters that result in the lowest number of bit transitions.

Note that this autoencoder network is decoupled from the regression network and design codes agnostic to classifiers’ characteristics for a given regression problem.

A.3 EXPERIMENTAL METHODOLOGY

We use 11 benchmarks covering four different regression tasks for evaluation. This section summarizes the experimental setup, including datasets, evaluation metrics, hyperparameters, and related work for each of these tasks. We also report the training time using an NVIDIA RTX 2080 Ti GPU with 11GB of memory for each benchmark.

A.3.1 HEAD POSE ESTIMATION

In landmark-free 2D head pose estimation, for a given 2D image, the head pose of a human is directly estimated in terms of three angles: yaw, pitch, and roll. We use loose cropping around the center with random flipping for data augmentation. We use the ResNet50 network as the feature extractor. This network is initialized using pre-trained parameters for ImageNet (Russakovsky et al., 2015) dataset. During the training for RLEL the entire network, including the feature extractor, is trained.

Datasets: We use the evaluation methodology followed by prior works (Ruiz et al., 2018; Yang et al., 2019). Two protocols are used for evaluation.

Protocol 1 (LFH1): This protocol uses the BIWI (Fanelli et al., 2013) dataset for training and evaluation. This dataset consists of 15, 128 frames of 20 subjects. Random 70% – 30% splits are used for training and evaluation. The ranges of yaw, pitch, and roll angles are $[-75^\circ, 75^\circ]$, $[-65^\circ, 85^\circ]$, and $[-55^\circ, 45^\circ]$, respectively.

Protocol 2 (LFH2): In this protocol, the network is trained using the 300W-LP (Zhu et al., 2016) dataset consisting of 122, 450 samples. AFLW2000 (Zhu et al., 2016) dataset is used for evaluation. The range of all labels is $[-99^\circ, 99^\circ]$ in this setting.

Evaluation metrics: We report the Mean Absolute Error (MAE) between the targets (y_i) and predictions (\hat{y}_i). Let N represent the number of samples, and P represent the number of labels (three

in head pose estimation). The MAE is defined as:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N \frac{1}{P} \sum_{j=1}^P |y_{i,j} - \hat{y}_{i,j}| \quad (12)$$

Network architecture and training parameters: Table 12 summarizes the hyperparameters used for RLEL. The learning rate of the decoding matrix D is kept $10\times$ higher than the learning rate of the feature extractor. L2 regularization with weight of 0.0001 is used for direct regression.

Table 12: Training parameters for LFH1.

Approach	Label range/ Quantization levels	Optimizer	Epochs	Batch size	Learning rate	Learning rate schedule	β	α	Training time (GPU hours)
LFH1	Yaw: $[-75^\circ, 75^\circ]/150$, Pitch: $[-65^\circ, 85^\circ]/150$, Roll: $[-55^\circ, 45^\circ]/100$	Adam, weight decay=0, momentum = 0	50	8	0.0001	1/10 after 30 Epochs	0.5	1.0	2
LFH2	$[-99^\circ, 99^\circ]/200$	Adam, weight decay=0, momentum = 0	20	16	0.00001	1/10 after 10 Epochs	2.0	0.0	4

Related work Head pose estimation is a widely studied problem. Existing task-specialized approaches propose different loss formulations or feature extractors to improve the error. HopeNet (Ruiz et al. 2018) proposed a combination of regression and classification loss. SSR-Net (Yang et al. 2018) and FSA-Net (Yang et al. 2019) proposed stage-wise soft regression. QuatNet (Hsu et al. 2019) proposed to use MSE loss with custom ordinal regression loss. RAFA-Net (Behera et al. 2021) proposed an attention-based feature extractor architecture. Table 13 and Table 14 compare the performance of RLEL with related work.

Table 13: Landmark-free 2D Head poses estimation evaluation for protocol 1 (HPE1 and HPE3).

Approach	Feature Extractor	#Params (M)	Yaw	Pitch	Roll	MAE
SSR-Net-MD (Yang et al. 2018) (Soft regression)	SSR-Net	1.1	4.24	4.35	4.19	4.26
FSA-Caps-Fusion (Yang et al. 2019) (Soft regression)	FSA-Net	5.1	2.89	4.29	3.60	3.60
RAFA-Net (Behera et al. 2021) (Direct Regression)	RAFA-Net	69.8	3.07	4.30	2.82	3.40
Direct regression (L2 loss)	ResNet50	23.5	3.80	4.63	4.28q	4.22 ± 0.35
BEL (Shah et al. 2022)	ResNet50	23.6	3.32	3.80	3.53	3.56 ± 0.11
RLEL	ResNet50	23.6	3.41	3.20	3.97	3.55 ± 0.10

Table 14: Landmark-free 2D Head poses estimation evaluation for protocol 2 (HPE2 and HPE4).

Approach	Feature Extractor	#Params (M)	Yaw	Pitch	Roll	MAE
SSR-Net-MD (Yang et al. 2018) (Soft regression)	SSR-Net	1.1	5.14	7.09	5.89	6.01
FSA-Caps-Fusion (Yang et al. 2019) (Soft regression)	FSA-Net	5.1	4.50	6.08	4.64	5.07
RAFA-Net (Behera et al. 2021) (Direct Regression)	RAFA-Net (HPE4)	69.8	3.60	4.92	3.88	4.13
HopeNet* ($\alpha = 2$) (Ruiz et al. 2018) (classification + regression loss)	ResNet50	23.9	6.47	6.56	5.44	6.16
Direct regression (L2 loss)	ResNet50	23.5	5.61	6.13	4.18	5.32 ± 0.12
BEL (Shah et al. 2022)	ResNet50	23.6	4.54	5.76	3.96	4.77 ± 0.05
RLEL	ResNet50	23.6	4.69	5.79	3.86	4.77 ± 0.05

A.3.2 FACIAL LANDMARK DETECTION

Facial landmark detection focuses on finding the (x, y) coordinates of facial keypoints for a given 2D image.

Evaluation metrics: We report the Normalized Mean Error (NME) between the targets y_i and predictions \hat{y}_i . Inter-ocular distance normalization is used for all datasets. For N test samples, P facial landmarks, and L normalization factor, the NME is defined as:

$$\text{NME} = \frac{1}{N} \sum_{i=1}^N \frac{1}{P} \cdot \frac{1}{L} \sum_{j=1}^P |y_{i,j} - \hat{y}_{i,j}|_2 \quad (13)$$

Datasets: We use three datasets widely used for facial landmark detection: COFW (Burgos-Artiz et al., 2013), 300W (Sagonas et al., 2013), and WFLW (Wu et al., 2018). HRNetV2-W18 network architecture for feature extraction (Wang et al., 2020) and the modified regressor architecture for label encoding proposed by BEL (Shah et al., 2022) are used in this work. Random flipping, scaling ($0.75 - 1.25$), and rotation (± 30) are used for data augmentation. The COFW dataset consists of 1,345 training and 507 testing images annotated with 29 landmarks. The training set of the 300W dataset has 3,148 images annotated with 68 facial landmarks. This dataset provides four test sets: full test set, common subset, challenging subset, and the official test set with 300 indoor and 300 outdoor images. WFLW dataset is a comparatively larger dataset with 7,500 training and 2,500 testing images. Each image is annotated with 98 facial landmarks. The test set is divided into six subsets: large pose, expression, illumination, make-up, occlusion, and blur.

Training parameters: Table 15 provides a summary of all the training parameters. The learning rate of the decoding matrix D is kept $20\times$ higher than the learning rate of the feature extractor. The HRNetV2-W18 network is initialized with pretrained weight parameters for the ImageNet dataset. We refer to HRNetV2-W18 evaluated on COFW as **FLD1/FLD1_s**, on 300W as **FLD2/FLD2_s**, and on WFLW as **FLD3/FLD3_s**.

Table 15: Training parameters for facial landmark detection for HRNetV2-W18 feature extractor.

Dataset/ Benchmark	Optimizer	Epochs	Batch size	Learning rate (BEL/Direct regres- sion/Multiclass classification)	Learning rate schedule	β	α	Training time (GPU hours)
COFW/ FLD1	Adam, weight decay=0, momentum = 0	60	8	0.0005/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	3.0	0.0	$\frac{1}{2}$
COFW/ FLD1_s	Adam, weight decay=0, momentum = 0	60	8	0.0005/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	4.0	0.0	$\frac{1}{2}$
300W/ FLD2	Adam, weight decay=0, momentum = 0	60	8	0.0007/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	5.0	1.0	3
300W/ FLD2_s	Adam, weight decay=0, momentum = 0	60	8	0.0007/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	5.0	0.05	3
WFLW/ FLD3	Adam, weight decay=0, momentum = 0	60	8	0.0003/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	0.0	0.1	5
WFLW/ FLD3_s	Adam, weight decay=0, momentum = 0	60	8	0.0003/0.0003/ 0.0003	1/10 after 30 and 50 Epochs	5.0	0.1	5

Related work Facial landmark detection is a widely studied problem. A common approach is to use heatmap regression, where the target heatmaps are generated by assuming a Gaussian distribution around the ground truth landmark location. Prior works proposed the use of binary heatmaps with pixel-wise binary cross-entropy loss (Bulat & Tzimiropoulos, 2016). HRNet (Wang et al., 2020) proposed a feature extractor that maintains high-resolution representations and uses heatmap regression. AWing (Wang et al., 2019) proposed a modified heatmap regression loss function with adapted wing loss. AnchorFace (Xu et al., 2020) used anchoring of facial landmarks on

templates. LUVLi (Kumar et al. 2020) proposed a landmark’s location, uncertainty, and visibility likelihood-based loss. Table 16-18 compare RLEL with related work.

Table 16: Facial landmark detection results on COFW dataset (FLD1).

Approach	Feature Extractor	#Params/ GFlops	Test NME	FR _{0.1}
LAB (w B) (Wu et al. 2018)	Hourglass	25.1/19.1	3.92	0.39
AWing (Wang et al. 2019)*	Hourglass	25.1/19.1	4.94	-
HRNetV2-W18 (Wang et al. 2020) (Heatmap regression)	HRNetV2-W18	9.6/4.6	3.45	0.19
Direct regression (L2 loss)	HRNetV2-W18	10.2/4.7	3.96 \pm 0.02	0.29
Direct regression (L1 loss)	HRNetV2-W18	10.2/4.7	3.60 \pm 0.02	0.29
BEL (Shah et al. 2022)	HRNetV2-W18	10.6/4.6	3.34 \pm 0.02	0.40
RLEL	HRNetV2-W18	10.6/4.6	3.36 \pm 0.01	0.20

*Uses different data augmentation for the training

Table 17: Facial landmark detection results on 300W dataset (FLD2).

Approach	Feature Extractor	#Params/ GFlops	Test	Common	Challenging	Full
DAN (Kowalski et al. 2017)	-	-	-	3.19	5.24	3.59
LAB (w B) (Wu et al. 2018)	Hourglass	25.1/19.1	-	2.98	5.19	3.49
AnchorFace (Xu et al. 2020)	ShuffleNet-V2	-	-	3.12	6.19	3.72
AWing (Wang et al. 2019)*	Hourglass	25.1/19.1	-	<u>2.72</u>	<u>4.52</u>	<u>3.07</u>
LUVLi (Kumar et al. 2020)	CU-Net	-	-	2.76	5.16	3.23
HRNetV2-W18 (Wang et al. 2020) (Heatmap regression)	HRNetV2-W18	9.6/4.6	-	2.87	5.15	3.32
Direct regression (L2 loss)	HRNetV2-W18	10.2/4.7	4.40	3.25	5.65	3.71 \pm 0.05
Direct regression (L1 loss)	HRNetV2-W18	10.2/4.7	4.26	3.10	5.42	3.54 \pm 0.03
BEL (Shah et al. 2022)	HRNetV2-W18	11.2/4.6	4.09	2.91	5.50	3.40 \pm 0.02
RLEL	HRNetV2-W18	11.2/4.6	4.03	2.90	5.39	3.37 \pm 0.02

*Uses different data augmentation for the training

Table 18: Facial landmark detection results (NME) on WFLW test (FLD3) and 6 subsets: pose, expression (expr.), illumination (illu.), make-up (mu.), occlusion (occu.) and blur. $\theta = 10$ is used for BEL.

Approach	Feature Extractor	#Params/ GFlops	Test	Pose	Expr.	Illu.	MU	Occu.	Blur
LAB (w B) (Wu et al. 2018)	Hourglass	25.1/19.1	5.27	10.24	5.51	5.23	5.15	6.79	6.32
AnchorFace (Xu et al. 2020)*	HRNetV2-W18	-/5.3	<u>4.32</u>	7.51	4.69	<u>4.20</u>	4.11	<u>4.98</u>	<u>4.82</u>
AWing (Wang et al. 2019)*	Hourglass	25.1/19.1	4.36	<u>7.38</u>	<u>4.58</u>	4.32	4.27	5.19	4.96
LUVLi (Kumar et al. 2020)	CU-Net	-	4.37	-	-	-	-	-	-
HRNetV2-W18 (Wang et al. 2020) (Heatmap regression)	HRNetV2-W18	9.6/4.6	4.60	7.94	4.85	4.55	4.29	5.44	5.42
Direct regression (L1 loss)	HRNetV2-W18	10.2/4.7	4.64 \pm 0.03	8.13	4.96	4.49	4.45	5.41	5.25
BEL (Shah et al. 2022)	HRNetV2-W18	11.7/4.6	4.36 \pm 0.02	7.53	4.64	4.28	4.19	5.19	5.05
RLEL	HRNetV2-W18	11.7/4.6	4.35 \pm 0.01	7.57	4.57	4.36	4.19	5.25	5.07

*Uses different data augmentation for the training

A.3.3 AGE ESTIMATION

This task focuses on predicting a person’s age from a given 2D image. MAE (Equation 12) and Cumulative Score (CS) are used as the evaluation metrics, and ResNet50 (He et al. 2016) is used as the feature extractor. CS θ is the percentage of test samples with absolute error less than θ years.

Datasets MORPH-II (Ricanek & Tesafaye, 2006) and AFAD (Niu et al. 2016) datasets are used for evaluation. We follow the protocols for preprocessing and data augmentation of datasets as per prior works (Shah et al. 2022; Raschka 2018). MORPH-II dataset consists of 55,608 images with random

split of 39,617 training, 4,398 validation, and 11,001 test images. The AFAD dataset consists of 164,432 images with random split of 118,492 training, 13,166 validation, and 32,763 test images.

Training parameters: Table 19 summarizes the training parameters for AE1 (MORPH-II) and AE2 (AFAD) benchmarks. The learning rate of the decoding matrix D is kept $10\times$ higher than the learning rate of the feature extractor. L2 regularization with weight of 0.001 is used for direct regression. Training for AE1 and AE2 consumes ~ 2 and ~ 7 hours, respectively.

Table 19: Training parameters for age estimation using MORPH-II and AFAD dataset

Bench-mark	Optimizer	Epochs	Batch size	Learning rate	Learning rate schedule	β	α
AE1	Adam, weight decay=0, momentum=0	50	64	0.0001	-	0.0	2.0
AE2	Adam, weight decay=0, momentum=0	50	64	0.0001	-	0.0	5.0

Related work Different approaches including ordinal regression (Niu et al., 2016; Cao et al., 2020; Pan et al., 2018; Gao et al., 2018), soft stage-wise regression (Yang et al., 2018; 2019), soft labels (Díaz & Marathe, 2019) have been proposed for age estimation. OR-CNN (Niu et al., 2016) and CORAL-CNN (Cao et al., 2020) proposed ordinal regression by binary classification with threshold-based encodings (i.e., unary codes). DLDL (Gao et al., 2018) augmented the loss function with KL-divergence between softmax output and soft target probability distributions. MV-Loss (Pan et al., 2018) proposed to penalize the prediction based on the variance of the age distribution. We compare CLL with related work in Table 20 and Table 21.

Table 20: Age estimation results on MORPH-II dataset (AE1).

Approach	Feature extractor	#Parameters (M)	MORPH-II (MAE)	MORPH-II (CS $\theta = 5$)
OR-CNN (Niu et al., 2016) (Ordinal regression by binary classification)	-	1.0	2.58	0.71
MV Loss (Pan et al., 2018) (Direct regression)	VGG-16	138.4	2.41	0.889
DLDL-v2 (Gao et al., 2018) (Ordinal regression with multi-class classification)	ThinAgeNet	3.7	1.96*	-
CORAL-CNN (Cao et al., 2020) (Ordinal regression by binary classification)	ResNet34	21.3	2.49	-
Direct Regression (L2 Loss)	ResNet50	23.1	2.37 \pm 0.01	0.903 \pm 0.002
BEL (Shah et al., 2022)	ResNet50	23.1	2.27 \pm 0.01	0.928 \pm 0.001
RLEL	ResNet50	23.1	2.28 \pm 0.01	0.901 \pm 0.002

Table 21: Age estimation results on AFAD dataset (AE2).

Approach	Feature extractor	#Parameters (M)	AFAD (MAE)	AFAD (CS $\theta = 5$)
OR-CNN (Niu et al., 2016) (Ordinal regression by binary classification)	-	1.0	3.51	0.74
CORAL-CNN (Cao et al., 2020) (Ordinal regression by binary classification)	ResNet34	21.3	3.47	-
Direct Regression (L2 Loss)	ResNet50	23.1	3.16 \pm 0.02	0.810 \pm 0.02
BEL (Shah et al., 2022)	ResNet50	23.1	3.11 \pm 0.01	0.823 \pm 0.001
RLEL	ResNet50	23.1	3.14 \pm 0.01	0.78 \pm 0.002

A.3.4 END-TO-END SELF DRIVING

For the regression task of end-to-end autonomous driving, we use the NVIDIA PilotNet dataset, and PilotNet model (Bojarski et al., 2016). In this task, for a given image of the road, the angle of the steering wheel that should be taken next is predicted. MAE (Equation 12) is used as the evaluation metric.

Dataset The PilotNet driving dataset consists of 45, 500 images taken around Rancho Palos Verdes and San Pedro, California (Chen). We use the data augmentation technique used by prior works (Bojarski et al., 2016).

Training parameters Table 22 summarizes the training parameters. The learning rate of the decoding matrix D is kept $10\times$ higher than the learning rate of the feature extractor.

Table 22: Training parameters for end-to-end autonomous driving using PilotNet.

Optimizer	Epochs	Batch size	Learning rate	Learning rate schedule	β	α
SGD with weight decay= $1e-5$, momentum=0	50	64	0.1	1/10 at 10, 30 epochs	0.0	2.0

Related work End-to-end autonomous driving is an interesting task with increasing attention. PilotNet (Bojarski et al., 2017) used a small, application-specific network. We compare RLEL with the baseline PilotNet architecture in Table 23.

Table 23: End-to-end autonomous driving results on PilotNet dataset (PN) and architecture (Bojarski et al., 2017).

Approach	Feature extractor	#Parameters (M)	MAE
PilotNet (Bojarski et al., 2017)	PilotNet	1.8	4.24 ± 0.45
BEL (Shah et al., 2022)	PilotNet	1.8	3.11 ± 0.01
RLEL	PilotNet	1.8	$\mathbf{2.94} \pm 0.01$

REFERENCES

- Erin L. Allwein, Robert E. Schapire, and Yoram Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *J. Mach. Learn. Res.*, 1:113–141, September 2001. doi: 10.1162/15324430152733133.
- Ardhendu Behera, Zachary Wharton, Pradeep Hewage, and Swagat Kumar. Rotation axis focused attention network (rafa-net) for estimating head pose. In *Computer Vision – ACCV 2020*, 2021.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to End Learning for Self-Driving Cars. *arXiv:1604.07316*, 2016.
- Mariusz Bojarski, Philip Yeres, Anna Choromanaska, Krzysztof Choromanski, Bernhard Firner, Lawrence Jackel, and Urs Muller. Explaining how a deep neural network trained with end-to-end learning steers a car. *arXiv:1704.07911*, 2017.
- Adrian Bulat and Georgios Tzimiropoulos. Human Pose Estimation via Convolutional Part Heatmap Regression. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision – ECCV 2016*, pp. 717–732, 2016.
- Xavier Burgos-Artizzu, Pietro Perona, and Piotr Dollár. Robust Face Landmark Estimation under Occlusion. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1513–1520, 12 2013. doi: 10.1109/ICCV.2013.191.
- Wenzhi Cao, Vahid Mirjalili, and Sebastian Raschka. Rank consistent ordinal regression for neural networks with application to age estimation. *Pattern Recognition Letters*, 140:325–331, 2020. doi: <https://doi.org/10.1016/j.patrec.2020.11.008>.
- Sully Chen. Driving-datasets. <https://github.com/SullyChen/driving-datasets>.
- M. Cissé, T. Artières, and Patrick Gallinari. Learning Compact Class Codes for Fast Inference in Large Multi Class Classification. In Peter A. Flach, Tijl De Bie, and Nello Cristianini (eds.), *Machine Learning and Knowledge Discovery in Databases*, pp. 506–520. Springer Berlin Heidelberg, 2012.

- Sajad Darabi, Mouloud Belbahri, Matthieu Courbariaux, and Vahid Partovi Nia. Regularized binary network training. *arXiv preprint arXiv:1812.11800*, 2018.
- T. G. Dietterich and G. Bakiri. Solving Multiclass Learning Problems via Error-Correcting Output Codes. *Journal of Artificial Intelligence Research*, 2:263–286, 1995. doi: 10.1613/jair.105.
- Raúl Díaz and Amit Marathe. Soft labels for ordinal regression. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. doi: 10.1109/CVPR.2019.00487.
- Stefan Falkner, Aaron Klein, and Frank Hutter. Bohb: Robust and efficient hyperparameter optimization at scale. In *ICML*, 2018.
- Gabriele Fanelli, Matthias Dantone, Juergen Gall, Andrea Fossati, and Luc Gool. Random forests for real time 3d face analysis. *International Journal of Computer Vision*, 101(3):437–458, February 2013.
- Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep Ordinal Regression Network for Monocular Depth Estimation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 2002–2011, 2018.
- A.E. Gamal, L. Hemachandra, I. Shperling, and V. Wei. Using simulated annealing to design good codes. *IEEE Transactions on Information Theory*, pp. 116–123, 1987. doi: 10.1109/TIT.1987.1057277.
- Bin-Bin Gao, Hong-Yu Zhou, Jianxin Wu, and Xin Geng. Age estimation using expectation of label distribution learning. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pp. 712–718, 7 2018. doi: 10.24963/ijcai.2018/99.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- Heng-Wei Hsu, Tung-Yu Wu, Sheng Wan, Wing Hung Wong, and Chen-Yi Lee. Quatnet: Quaternion-based head pose estimation with multiregression loss. *IEEE Transactions on Multimedia*, 21(4): 1035–1046, 2019. doi: 10.1109/TMM.2018.2866770.
- Lu Jin, Xiangbo Shu, Kai Li, Zechao Li, Guo-Jun Qi, and Jinhui Tang. Deep ordinal hashing with spatial attention. *Transaction on Image Processing*, 28(5):2173–2186, 2019.
- Marek Kowalski, Jacek Naruniec, and Tomasz Trzcinski. Deep alignment network: A convolutional neural network for robust face alignment. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 2034–2043, 2017. doi: 10.1109/CVPRW.2017.254.
- Abhinav Kumar, Tim K. Marks, Wenxuan Mou, Ye Wang, Michael Jones, Anoop Cherian, Toshiaki Koike-Akino, Xiaoming Liu, and Chen Feng. LUVLi face alignment: Estimating Landmarks’ location, uncertainty, and visibility likelihood. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2020. doi: 10.1109/CVPR42600.2020.00826.
- H. Lai, Y. Pan, Ye Liu, and S. Yan. Simultaneous feature learning and hash coding with deep neural networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- Ling Li and Hsuan-Tien Lin. Ordinal regression by extended binary classification. In *Proceedings of the 19th International Conference on Neural Information Processing Systems*, pp. 865–872, 2006.
- Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *J. Mach. Learn. Res.*, 18(1): 6765–6816, jan 2017. ISSN 1532-4435.
- William Libaw and Leonard Craig. A photoelectric decimal-coded shaft digitizer. *Electronic Computers, Transactions of the I.R.E. Professional Group on*, EC-2:1 – 4, 10 1953.
- Haomiao Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Deep supervised hashing for fast image retrieval. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

- Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020.
- Xiao Luo, Haixin Wang, Daqing Wu, Chong Chen, Minghua Deng, Jianqiang Huang, and Xian-Sheng Hua. A survey on deep hashing methods. *ACM Trans. Knowl. Discov. Data*, 2022. ISSN 1556-4681. doi: 10.1145/3532624.
- Zhenxing Niu, Mo Zhou, Le Wang, Xinbo Gao, and Gang Hua. Ordinal regression with multiple output CNN for age estimation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016.
- Mohammad Norouzi, David J. Fleet, and Ruslan Salakhutdinov. Hamming distance metric learning. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, 2012.
- Hongyu Pan, Hu Han, Shiguang Shan, and Xilin Chen. Mean-variance loss for deep age estimation from a face. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5285–5294, 2018. doi: 10.1109/CVPR.2018.00554.
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 5301–5310, 2019.
- Sebastian Raschka. MLxtend: Providing machine learning and data science utilities and extensions to Python’s scientific computing stack. *Journal of Open Source Software*, 3(24):638, April 2018. doi: 10.21105/joss.00638.
- K. Ricanek and T. Tesafaye. Morph: a longitudinal image database of normal adult age-progression. In *7th International Conference on Automatic Face and Gesture Recognition (FGR06)*, pp. 341–345, 2006. doi: 10.1109/FGR.2006.78.
- Nataniel Ruiz, Eunji Chong, and James M. Rehg. Fine-grained head pose estimation without key-points. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.
- Olga Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Zhiheng Huang, A. Karpathy, A. Khosla, M. Bernstein, A. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211–252, 2015.
- C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *2013 IEEE International Conference on Computer Vision Workshops*, pp. 397–403, 2013.
- Deval Shah, Zi Yu Xue, and Tor M. Aamodt. Label encoding for regression networks. In *International Conference on Learning Representations*, April 2022. URL <https://openreview.net/pdf?id=8WawVDdKq1L>.
- Yang Song, Qiyu Kang, and Wee Peng Tay. Error-Correcting Output Codes with Ensemble Diversity for Robust Learning in Neural Networks. *AAAI*, 2021.
- Gunjan Verma and Ananthram Swami. Error correcting output codes improve probability estimation and adversarial robustness of deep neural networks. *Advances in Neural Information Processing Systems*, 32(NeurIPS), 2019.
- J. Wang, T. Zhang, J. Song, N. Sebe, and H. Shen. A survey on learning to hash. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(04), 2018.
- Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Minghui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, PP, April 2020.

- Xinyao Wang, Liefeng Bo, and Fuxin Li. Adaptive wing loss for robust face alignment via heatmap regression. In *2019 IEEE International Conference on Computer Vision (ICCV)*, pp. 6970–6980, 2019. doi: 10.1109/ICCV.2019.00707.
- W. Wu, Chen Qian, S. Yang, Q. Wang, Y. Cai, and Qiang Zhou. Look at boundary: A boundary-aware face alignment algorithm. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2129–2138, 2018.
- Rongkai Xia, Yan Pan, Hanjiang Lai, Cong Liu, and Shuicheng Yan. Supervised hashing for image retrieval via image representation learning. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, AAAI’14, 2014.
- Zixuan Xu, Banghuai Li, Miao Geng, Ye Yuan, and Gang Yu. Anchorface: An anchor-based facial landmark detector across large poses. *ArXiv:2007.03221*, 2020.
- Tsun-Yi Yang, Yi-Hsuan Huang, Yen-Yu Lin, Pi-Cheng Hsiu, and Yung-Yu Chuang. Ssr-net: A compact soft stagewise regression network for age estimation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, pp. 1078–1084. AAAI Press, 2018. ISBN 9780999241127.
- Tsun-Yi Yang, Yi-Ting Chen, Yen-Yu Lin, and Yung-Yu Chuang. Fsa-net: Learning fine-grained structure aggregation for head pose estimation from a single image. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4578–4587, 2021.
- Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and Stan Li. Face alignment across large poses: A 3d solution. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 146–155, 06 2016.