
Model Distillation with Knowledge Transfer from Face Classification to Alignment and Verification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Knowledge distillation is a potential solution for model compression. The idea
2 is to make a small student network imitate the target of a large teacher network,
3 then the student network can be competitive to the teacher one. Most previous
4 studies focus on model distillation in the classification task, where they propose
5 different architectures and initializations for the student network. However, only the
6 classification task is not enough, and other related tasks such as regression and
7 retrieval are barely considered. To solve the problem, in this paper, we take face
8 recognition as a breaking point and propose model distillation with knowledge
9 transfer from face classification to alignment and verification. By selecting appropriate
10 initializations and targets in the knowledge transfer, the distillation can be
11 easier in non-classification tasks. Experiments on the CelebA and CASIA-WebFace
12 datasets demonstrate that the student network can be competitive to the teacher
13 one in alignment and verification, and even surpasses the teacher network under
14 specific compression rates. In addition, to achieve stronger knowledge transfer, we
15 also use a common initialization trick to improve the distillation performance of
16 classification. Evaluations on the CASIA-Webface and large-scale MS-Celeb-1M
17 datasets show the effectiveness of this simple trick.

18 1 Introduction

19 Since the emergence of Alexnet[19], larger and deeper networks have shown to be more powerful[31,
20 32, 13]. However, as the network going larger and deeper, it becomes difficult to use it in mobile
21 devices. Therefore, model compression has become necessary in compressing the large network into
22 a small one. In recent years, many compression methods have been proposed, including knowledge
23 distillation[24, 3, 1, 35, 22, 29, 14, 28], weight quantization[8, 6, 26], weight pruning[12, 10, 11, 16,
24 33] and weight decomposition[40, 4, 7, 34, 23]. In this paper, we focus on the knowledge distillation,
25 which is a potential approach for model compression.

26 In knowledge distillation, there is usually a large teacher network and a small student one, and the
27 objective is to make the student network competitive to the teacher one by learning the specific target
28 of the teacher network. Previous studies mainly consider the selection of targets in the classification
29 task, *e.g.*, hidden layers[22], logits[1, 35, 29] or soft predictions[14, 28]. However, only the distillation
30 of the classification task is not enough, and some common tasks such as regression and retrieval
31 should also be considered. In this paper, we take face recognition as a breaking point that we start with
32 the knowledge distillation in face classification, and consider the distillation on two domain-similar
33 tasks, including face alignment and verification. The objective of face alignment is to locate the
34 key-point locations in each image; while in face verification, we have to determine if two images
35 belong to the same identity.

36 For distillation on non-classification tasks, one intuitive idea is to adopt a similar method as in face
37 classification that trains teacher and student networks from scratch. In this way, the distillation on all
38 these tasks will be independent, and this is a possible solution. However, this independence cannot
39 give the best distillation performance. There has been strong evidence that in object detection[27],
40 object segmentation[5] and image retrieval[41], they all used the pretrained classification model(on
41 ImageNet) as initialization to boost performance. This success comes from the fact that their
42 domains are similar, which makes them transfer a lot from low-level to high-level representation[39].
43 Similarly, face classification, alignment and verification also share the similar domain, thus we
44 transfer the distilled knowledge of classification by taking its teacher and student networks to
45 initialize corresponding networks in alignment and verification.

46 Another problem in knowledge transfer is what targets should be used for distillation? In face
47 classification, the knowledge is distilled from the teacher network by learning its soft-prediction,
48 which has been proved to work well[14, 28]. However, in face alignment[37] and verification[37],
49 they have additional task-specific targets for learning. As a result, selecting the classification or
50 task-specific target for distillation remains a problem. One intuitive idea is to measure the relevance
51 of objectives between non-classification and classification tasks. For example, it is not obvious to
52 see the relation between face classification and alignment, but the classification can help a lot in
53 verification. Therefore, it seems reasonable that if the tasks are highly related, the classification target
54 is preferred, or the task-specific target is better.

55 Inspired by the above thoughts, in this paper, we propose the model distillation in face alignment
56 and verification by transferring the distilled knowledge from face classification. With appropriate
57 selection of initializations and targets, we show that the distillation performance of alignment and
58 verification on the CelebA[21] and CASIA-WebFace[38] datasets can be largely improved, and the
59 student network can even exceed the teacher network under specific compression rates.

60 This knowledge transfer is our main contribution. In addition, we realize that in the proposed
61 method, the knowledge transfer depends on the distillation of classification, thus we use a common
62 initialization trick to further boost the distillation performance of classification at the beginning.
63 Evaluations on the CASIA-WebFace[38] and large-scale MS-Celeb-1M[9] datasets show that this
64 simple trick can give the best distillation results in the classification task.

65 2 Related Work

66 In this part, we introduce some previous studies on knowledge distillation. Particularly, all the
67 following studies focus on the task of classification. Buciluă *et al.*[3] propose to generate synthetic
68 data by a teacher network, then a student network is trained with the synthetic data to mimic the
69 identity labels of the teacher network. However, Ba and Caruana[1] observe that these labels have
70 lost the uncertainties of the teacher network, thus they propose to regress the logits (pre-softmax
71 activations)[14]. Besides, they prefer the student network to be deep, which is good to mimic complex
72 functions. To better learn the function, Gregor *et al.*[35] observe the student network should not
73 only be deep, but also be convolutional, and they get competitive performance to the teacher network
74 in CIFAR[18]. Most methods need a large ensemble of teacher networks for distillation, but this
75 will take a long training and inference time[29]. To address the issue, Sau and Balasubramanian[29]
76 propose noise-based regularization that can simulate the logits of multiple teacher networks. However,
77 Luo *et al.*[22] observe the values of the logits are unconstrained, and the high dimensionality will
78 also cause fitting problem. As a result, they use the hidden layer as it captures as much information
79 as the logits but is more compact.

80 All these methods only use the targets of the teacher network in distillation, while if the target is not
81 confident, the training of the student network will be difficult. To solve the problem, Hinton *et al.*[14]
82 propose a multi-task approach which uses identity labels and the target of the teacher network jointly.
83 Particularly, they use the post-softmax activations with temperature smoothing as the target, which
84 can better represent the label distribution. One problem is that student networks are mostly trained
85 from scratch. Given the fact that initialization is important, Romero *et al.*[28] propose to initialize
86 the shallow layers of the student network by regressing the mid-level target of the teacher network.
87 However, these studies only consider the knowledge distillation in classification, which largely limits
88 its application in model compression. In this paper, we consider the face recognition problem as a
89 breaking point and extend the distillation to non-classification tasks.

90 3 Distillation of Classification

91 Due to the proposed knowledge transfer depends on the distillation in classification, improving the
92 classification itself is necessary. In this part, we first review the idea of distillation for classification,
93 then introduce how to boost it by a simple initialization trick.

94 3.1 Review of Knowledge Distillation

95 We adopt the distillation framework in Hinton *et al.*[14], which is summarized as follows. Let T and S
96 be the teacher and student network, and their post-softmax predictions to be $\mathbf{P}_T = \text{softmax}(\mathbf{a}_T)$ and
97 $\mathbf{P}_S = \text{softmax}(\mathbf{a}_S)$, where \mathbf{a}_T and \mathbf{a}_S are the pre-softmax predictions, also called the logits[1, 35, 29].
98 However, the post-softmax predictions have lost some relative uncertainties that are more informative,
99 thus a temperature parameter τ is used to smooth predictions \mathbf{P}_T and \mathbf{P}_S to be \mathbf{P}_T^τ and \mathbf{P}_S^τ , which
100 are denoted as *soft predictions*:

$$\mathbf{P}_T^\tau = \text{softmax}(\mathbf{a}_T/\tau), \quad \mathbf{P}_S^\tau = \text{softmax}(\mathbf{a}_S/\tau). \quad (1)$$

101 Then, consider \mathbf{P}_T^τ as the target, knowledge distillation optimizes the following loss function

$$L(\mathbf{W}_S^{\text{cls}}) = H(\mathbf{P}_S, \mathbf{y}^{\text{cls}}) + \alpha H(\mathbf{P}_S^\tau, \mathbf{P}_T^\tau), \quad (2)$$

102 wherein $\mathbf{W}_S^{\text{cls}}$ is the parameter of the student network, and \mathbf{y}^{cls} is the identity label. For simplicity,
103 we omit *min* and the number of samples N , and denote the upper right symbol *cls* as the classification
104 task. In addition, $H(\cdot, \cdot)$ is the cross-entropy, thus the first term is the softmax loss, while the second
105 one is the cross-entropy between the soft predictions of the teacher and student network, with α
106 balancing between the two terms. This multi-task training is advantageous because the target \mathbf{P}_T^τ
107 cannot be guaranteed to be always correct, and if the target is not confident, the identity label \mathbf{y}^{cls}
108 will take over the training of the student work.

109 3.2 Initialization Trick

110 It is noticed that in Eqn.(2), the student network is trained from scratch. As demonstrated in [1, 35] that
111 deeper student networks are better for distillation, initialization thus has become very important[15,
112 2, 17]. Based on the evidence, Fitnet[28] first initializes the shallow layers of the student network
113 by regressing the mid-level target of the teacher network, then it follows Eqn.(2) for distillation.
114 However, only initializing the shallow layers is still difficult to learn high-level representation, which
115 is generated by deep layers. Furthermore, [39] shows that the network transferability increases as
116 tasks become more similar. In our case, the initialization and distillation are both classification tasks
117 with exactly the same data and identity labels, thus more deep layers should be initialized for higher
118 transferability, and we use a simple trick to achieve this.

119 To obtain an initial student network, we first train it with softmax loss:

$$L(\mathbf{W}_{S_0}^{\text{cls}}) = H(\mathbf{P}_S, \mathbf{y}^{\text{cls}}), \quad (3)$$

120 wherein the lower right symbol S_0 denotes the initialization for student network S . In this way, the
121 student network is fully initialized. Then, we modify Eqn.(2) as

$$L(\mathbf{W}_S^{\text{cls}} | \mathbf{W}_{S_0}^{\text{cls}}) = H(\mathbf{P}_S, \mathbf{y}^{\text{cls}}) + \alpha H(\mathbf{P}_S^\tau, \mathbf{P}_T^\tau), \quad (4)$$

122 wherein $\mathbf{W}_S^{\text{cls}} | \mathbf{W}_{S_0}^{\text{cls}}$ indicates that $\mathbf{W}_S^{\text{cls}}$ is trained with the initialization of $\mathbf{W}_{S_0}^{\text{cls}}$, and the two
123 entropy terms remain the same. This process is shown in Fig.1(a). It can be seen that the only
124 difference with Eqn.(2) is that the student network is trained with the full initialization, and this
125 simple trick has been commonly used, *e.g.*, initializing the VGG-16 model[31]. We later show that
126 this trick can get promising improvements over Eqn.(2) and Fitnet[28].

127 4 Distillation Transfer

128 In this part, we show how to transfer the distilled knowledge from face classification to face alignment
129 and verification. The knowledge transfer consists of two steps: transfer initialization and target
130 selection, which are elaborated as follows.

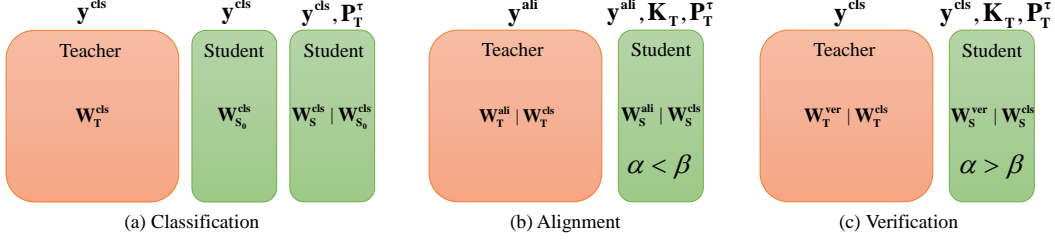


Figure 1: The pipeline of knowledge distillation in face classification, alignment and verification. \mathbf{y}^{cls} and \mathbf{y}^{ali} are the ground truth labels for classification and alignment respectively.

131 4.1 Transfer Initialization

132 The first step of transfer is the initialization. The motivation is based on the evidence that in detection,
 133 segmentation and retrieval, they have used the pretrained classification model (on ImageNet) as
 134 initialization to boost performance[27, 5, 41]. The availability of this idea comes from the fact
 135 that they share the similar domain, which makes them transfer easily from low-level to high-level
 136 representation[39]. Similarly, the domains of face classification, alignment and verification are also
 137 similar, thus we can transfer the distilled knowledge in the same way.

138 For simplicity, we denote the parameters of teacher and student networks in face classification as
 139 $\mathbf{W}_T^{\text{cls}}$ and $\mathbf{W}_S^{\text{cls}}$. Analogously, they are $\mathbf{W}_T^{\text{ali}}$ and $\mathbf{W}_S^{\text{ali}}$ in alignment, while $\mathbf{W}_T^{\text{ver}}$ and $\mathbf{W}_S^{\text{ver}}$ in
 140 verification. As analyzed above, in distillation of alignment and verification, teacher and student
 141 networks will be initialized by $\mathbf{W}_T^{\text{cls}}$ and $\mathbf{W}_S^{\text{cls}}$ respectively.

142 4.2 Target Selection

143 Based on the initialization, the second step is to select appropriate targets in the teacher network for
 144 distillation. One problem is that non-classification tasks have their own task-specific targets, but
 145 given the additional soft predictions \mathbf{P}_T^τ , which one should we use? To be clear, we first propose the
 146 general distillation for non-classification tasks as follows:

$$L(\mathbf{W}_S | \mathbf{W}_S^{\text{cls}}) = \Phi(\mathbf{W}_S, \mathbf{y}) + \alpha H(\mathbf{P}_S^\tau, \mathbf{P}_T^\tau) + \beta \Psi(\mathbf{K}_S, \mathbf{K}_T), \quad (5)$$

147 where \mathbf{W}_S and \mathbf{y} denote the task-specific network parameter and label respectively. $\Phi(\mathbf{W}_S, \mathbf{y})$ is
 148 the task-specific loss function, and $\Psi(\mathbf{K}_S, \mathbf{K}_T)$ is the task-specific distillation term with the targets
 149 selected as \mathbf{K}_T and \mathbf{K}_S in teacher and student networks. Besides, α and β are the balancing terms
 150 between classification and non-classification tasks. In Eqn.(5), the above problem has become how
 151 to set α and β for a given non-classification task. In the following two parts, we will give some
 152 discussions on two tasks: face alignment and verification.

153 4.2.1 Alignment

154 The task of face alignment is to locate the key-point locations for each image. Without loss of
 155 generality, there is no any identity labels, but only the keypoint locations for each image. Face
 156 alignment is usually considered as a regression problem[37], thus we train the teacher network with
 157 optimizing the Euclidean loss:

$$L(\mathbf{W}_T^{\text{ali}} | \mathbf{W}_T^{\text{cls}}) = \|\mathbf{R}_T - \mathbf{y}^{\text{ali}}\|^2, \quad (6)$$

158 wherein \mathbf{R}_T is the regression prediction of the teacher network and \mathbf{y}^{ali} is the regression label. In
 159 distillation, except for the available soft predictions \mathbf{P}_T^τ (classification target), another one is the
 160 task-specific target that can be the hidden layer \mathbf{K}_T [22], and it satisfies $\mathbf{R}_T = fc(\mathbf{K}_T)$ with fc
 161 being a fully-connected mapping.

162 In face classification, the key in distinguishing different identities is the appearance around the
 163 key-points such as shape and color, but the difference of key-point locations for different identities
 164 is tiny. As a result, face identity is not the main influencing factor for these locations, but it is still
 165 related as different identities may have slightly different locations. Instead, pose and viewpoint

166 variations have a much larger influence. Therefore, in face alignment, the hidden layer is preferred
 167 for distillation, which gives Eqn.(7) by setting $\alpha < \beta$, as shown in Fig.1(b).

$$L(\mathbf{W}_S^{\text{ali}}|\mathbf{W}_S^{\text{cls}}) = \|\mathbf{R}_S - \mathbf{y}^{\text{ali}}\|^2 + \alpha H(\mathbf{P}_S^\tau, \mathbf{P}_T^\tau) + \beta \|\mathbf{K}_S - \mathbf{K}_T\|^2. \quad (7)$$

168 4.2.2 Verification

169 The task of face verification is to determine if two images belong to the same identity. In verification,
 170 triplet loss[25, 30] is a widely used metric learning method[30], and we take it for model distillation.
 171 Without loss of generality, we have the same identity labels as in classification[38, 9], then the teacher
 172 network can be trained as

$$L(\mathbf{W}_T^{\text{ver}}|\mathbf{W}_T^{\text{cls}}) = \left[\|\mathbf{K}_T^a - \mathbf{K}_T^p\|^2 - \|\mathbf{K}_T^a - \mathbf{K}_T^n\|^2 + \lambda \right]_+, \quad (8)$$

173 where \mathbf{K}_T^a , \mathbf{K}_T^p and \mathbf{K}_T^n are the hidden layers for the anchor, positive and negative samples respec-
 174 tively, *i.e.*, a and p have the same identity, while a and n come from different identities. Besides, λ
 175 controls the margin between positive and negative pairs.

176 Similar to face alignment, we consider the hidden layer \mathbf{K}_T and soft prediction \mathbf{P}_T^τ as two possible
 177 targets in distillation. In fact, classification focuses on the difference of identities, *i.e.* the inter-class
 178 relation, and this relation can easily tell if two image have the same identity. As a result, classification
 179 can be beneficial to boost the performance of verification. Therefore, in face verification, the soft
 180 predictions are preferred for distillation, which gives the following loss function by setting $\alpha > \beta$, as
 181 shown in Fig.1(c).

$$L(\mathbf{W}_S^{\text{ver}}|\mathbf{W}_S^{\text{cls}}) = \left[\|\mathbf{K}_S^a - \mathbf{K}_S^p\|^2 - \|\mathbf{K}_S^a - \mathbf{K}_S^n\|^2 + \lambda \right]_+ + \alpha H(\mathbf{P}_S^\tau, \mathbf{P}_T^\tau) + \beta \|\mathbf{K}_S - \mathbf{K}_T\|^2. \quad (9)$$

182 Particularly, some studies[36] also shows the benefits by using softmax loss in Eqn.(8). For compari-
 183 son, we also add the softmax loss $H(\mathbf{P}_T, \mathbf{y}^{\text{cls}})$ and $H(\mathbf{P}_S, \mathbf{y}^{\text{cls}})$ in Eqn.(8) and Eqn.(9) respectively
 184 for further enhancement.

185 4.2.3 A Short Summary

186 As analyzed above, α and β should be set differently in the distillation of different tasks. The key is
 187 to measure the relevance of objectives between classification and non-classification tasks. For a given
 188 task, if the classification is highly related, then $\alpha > \beta$ is necessary, or $\alpha < \beta$ should be set. Though
 189 this rule cannot be theoretically guaranteed, it provides some guidelines to use knowledge distillation
 190 in more non-classification tasks.

191 5 Experimental Evaluation

192 In this section, we give the experimental evaluation of the proposed method. We first introduce the
 193 experimental setup in detail, and then show the results of knowledge distillation in the tasks of face
 194 classification, alignment and verification.

195 5.1 Experimental Setup

196 **Database:** We use three popular datasets for evaluation, including CASIA-WebFace[38], CelebA[21]
 197 and MS-Celeb-1M[9]. CASIA-WebFace contains 10575 people and 494414 images, while CelebA
 198 has 10177 people with 202599 images and the label of 5 key-point locations. Compared to the
 199 previous two, MS-Celeb-1M is a large-scale dataset that contains 100K people with 8.4 million
 200 images. In experiments, we use CASIA-WebFace and MS-Celeb-1M for classification, CelebA for
 201 alignment and CASIA-WebFace for verification.

202 **Evaluation:** In all datasets, we randomly split them into 80% training and 20% testing samples. In
 203 classification, we evaluate the top1 accuracy based on if the identity of the maximum prediction
 204 matches the correct identity label[19], and the results on the LFW[20] database (6000 pairs) are
 205 also reported by computing the percentage of how many pairs are correctly verified. In alignment,
 206 the Normalized Root Mean Squared Error (NRMSE) is widely used to evaluate the alignment
 207 performance[37]; while in verification, we compute the Euclidean distance between each pair in

208 testing samples, and the top1 accuracy is reported based on if a test sample and its nearest sample
 209 belong to the same identity. Particularly, LFW is not used in verification because 6000 pairs are not
 210 enough to see the difference obviously for different methods.

211 **Teacher and Student:** To learn the large number of identities, we use ResNet-50[13] as the teacher
 212 network, which is deep enough to handle our problem. For student networks, given the fact that deep
 213 student networks are better for knowledge distillation[1, 35, 28], we remain the same depth but divide
 214 the number of convolution kernels in each layer by 2, 4 and 8, which give ResNet-50/2, ResNet-50/4
 215 and ResNet-50/8 respectively.

216 **Pre-processing and Training:** Given an image, we resize it to 256×256 wherein a sub-image with
 217 224×224 is randomly cropped and flipped. Particularly, we use no mean subtraction or image
 218 whitening, as we use batch normalization right after the input data. In training, the batchsize is set
 219 to be 256, 64 and 128 for classification, alignment and verification respectively, and the Nesterov
 220 Accelerated Gradient(NAG) is adopted for faster convergence. For the learning rate, if the network is
 221 trained from scratch, 0.1 is used; while if the network is initialized, 0.01 is used to continue, and 30
 222 epochs are used in each rate. Besides, in distillation, student networks are trained with the targets
 223 of the teacher network generated online, and the temperature τ and margin λ are set to be 3 and 0.4
 224 by cross-validation. Finally, the balancing terms α and β have many possible combinations, and we
 225 show later how to avoid this by an experimental trick.

226 **Symbols in Experiments:** (1)*Scratch*: student networks are not initialized; (2)*Pretrain*: student
 227 networks are trained with task-specific initialization; (3)*Distill*: student networks are initialized with
 228 $\mathbf{W}_S^{\text{cls}}$; (4)*Soft*: the soft prediction \mathbf{P}_T^τ ; (5)*Hidden*: the hidden layer \mathbf{K}_T .

229 5.2 Comparison to Previous Studies

230 In this part, we compare the initialization trick to previous studies in classification. Table.1 shows the
 231 comparison of different targets and initializations. It can be observed from the first table that without
 232 any initialization, soft predictions achieve the best distillation performance, *i.e.*, 61.27%. Based on
 233 the best target, the second table gives the results of different initializations in distillation. We see that
 234 our full initialization obtains the best accuracy of 75.06%, which is much higher than other methods,
 235 *i.e.*, 10% and 5% higher than the *Scratch* and Fitnet[28]. This shows that the full initialization
 236 of student networks can give the highest transferability in classification, and also demonstrates the
 effectiveness of this simple trick.

Table 1: The comparison to previous studies with different initializations and targets. Results are given on CASIA-WebFace.

<i>CASIA-WebFace</i> Teacher Network		Student Network with Different Targets				
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	88.61	ResNet-50/8	Scratch	Hidden[22]	Logits[29]	Soft[14]

<i>CASIA-WebFace</i> Teacher Network		Student Network with Different Initializations				
Top1 acc(%)	ResNet-50	Network	Learning	Initialization		
	88.61	ResNet-50/8	Distill	Scratch[14]	Fitnet[28]	Ours

237

238 5.3 Face Classification

239 Base on the best initialization and target, Table.2 shows the distillation results of face classification
 240 on CASIA-WebFace and MS-Celeb-1M, and we have three main observations. Firstly, the student
 241 networks trained with full initialization can obtain large improvements over the ones trained from
 242 scratch, which further demonstrates the effectiveness of the initialization trick in large scale cases.
 243 Secondly, some student networks can be competitive to the teacher network or even exceed the teacher
 244 one by a large margin, *e.g.*, in CASIA-WebFace, ResNet-50/4 can be competitive to the teacher
 245 network, while ResNet-50/2 is about 3% higher than the teacher one in both top1 and LFW accuracy.
 246 Finally, in the large-scale MS-Celeb-1M, student networks cannot exceed the teacher network but

247 only be competitive, which shows that the knowledge distillation is still challenging in distilling a
 248 large number of identities.

Table 2: The top1 and LFW accuracy of knowledge distillation in classification. Results are obtained on CASIA-WebFace and MS-Celeb-1M.

<i>CASIA-WebFace</i>	Teacher Network		Student Network			
	ResNet-50		Initialization	ResNet-50/2	ResNet-50/4	ResNet-50/8
Top1 acc(%)	88.61		Scratch[14]	82.25	79.36	66.12
			Ours	91.01	87.21	75.06
LFW acc(%)	97.67		Scratch[14]	97.27	96.7	95.12
			Ours	98.2	97.57	96.18

<i>MS-Celeb-1M</i>	Teacher Network		Student Network			
	ResNet-50		Initialization	ResNet-50/2	ResNet-50/4	ResNet-50/8
Top1 acc(%)	90.53		Scratch[14]	84.59	81.94	57.84
			Ours	88.38	85.26	70.98
LFW acc(%)	99.11		Scratch[14]	98.61	98.03	96.33
			Ours	98.88	98.18	96.98

248

249 5.4 Face Alignment

250 In this part, we give the evaluation of distillation in face alignment. Table.3 shows the distillation
 251 results of ResNet-50/8 with different initializations and targets on CelebA. The reason we only
 252 consider ResNet-50/8 is that face alignment is a relatively easy problem and most studies use shallow
 253 and small networks, thus a large compression rate is necessary for the deep ResNet-50. One important
 254 thing is how to set α and β in Eqn.(7). As there are many possible combinations, we use a simple
 255 trick by measuring their individual influence and discard the target with the negative impact by setting
 256 $\alpha = 0$ or $\beta = 0$; while if they both have positive impacts, $\alpha > 0, \beta > 0$ should be set to keep both
 257 targets in distillation.

258 As shown in Table.3, when the initializations of *Pretrain* and *Distill* are used, $\alpha = 1, \beta = 0$ (soft
 259 prediction) always decreases performance, while $\alpha = 0, \beta = 1$ (hidden layer) gets consistent
 260 improvements, which implies that the hidden layer($\alpha = 0, \beta = 1$) is preferred in the distillation of
 261 face alignment. Therefore, $\alpha = 0, \beta = 1$ is used in face alignment. It can be observed in Table.3
 262 that *Distill* has a lower error rate than *Pretrain*, which shows that $\mathbf{W}_S^{\text{cls}}$ has higher transferability
 263 on high-level representation than the task-specific initialization. Besides, the highest distillation
 264 performance 3.21% is obtained by ResNet-50/8 with *Distill* and $\alpha = 0, \beta = 1$, and it can be
 competitive to the one of the teacher network(3.02%).

Table 3: The NRMSE(%) of knowledge distillation in face alignment with different initializations and targets. Results are obtained on CelebA.

<i>CelebA</i>	Teacher Network		Student Network				
	ResNet-50		Network	Initialization	Targets		
					$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$	$\alpha = 1, \beta = 0$
NRMSE(%)	3.02		ResNet-50/8	Pretrain	3.36	3.24	3.60
				Distill	3.29	3.21	3.54

265

266 5.5 Face Verification

267 In this part, we give the evaluation of distillation in face verification. Similar to alignment, we select
 268 α and β in the same way. Table.4 shows the verification results of different initializations and targets
 269 on CASIA-WebFace, and the results are given by Eqn.(9). It can be observed that no matter which
 270 student network or initialization is used, $\alpha = 0, \beta = 1$ (hidden layer) always decreases the baseline
 271 performance, while $\alpha = 1, \beta = 0$ (soft prediction) remains almost the same. As a result, we discard
 272 the hidden layer and only use the soft prediction.

273 One interesting observation in Table.4 is that $\alpha = 0, \beta = 0$ always get the best performance, and
 274 the targets do not work at all. One possible reason is that the target in classification is not confident,
 275 *i.e.*, the top1 accuracy of ResNet-50 in classification is only 88.61%. To improve the classification
 276 ability, we add additional softmax loss in Eqn.(8) and Eqn.(9), and the results are shown in Table.5.
 277 We see that the accuracy of ResNet-50/2 and ResNet-50/4 has obtained remarkable improvements,
 278 which implies that the classification targets that are not confident cannot help the distillation. But
 279 with the additional softmax loss, the student work can adjust the learning by identity labels. As a
 280 result, $\alpha = 1, \beta = 0$ can get the best distillation performance, which is even much higher than the
 teacher network, *e.g.*, 79.96% of ResNet-50/2 with *Distill* and $\alpha = 1, \beta = 0$.

Table 4: The top1 accuracy of distillation with single triplet loss in verification. Results are obtained on CASIA-WebFace.

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	73.81			ResNet-50/2	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	63.98	60.66	66.50
			Distill	71.29	68.74	71.23

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	73.81			ResNet-50/4	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	61.74	61.71	62.64
			Distill	68.17	66.74	68.12

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	73.81			ResNet-50/8	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	51.03	49.19	51.76
			Distill	56.69	53.99	56.52

281 Table 5: The top1 accuracy of distillation with joint triplet loss and softmax loss in verification. Results are obtained on CASIA-WebFace.

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	74.16			ResNet-50/2	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	72.38	70.54	73.62
			Distill	79.51	77.63	79.96

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	74.16			ResNet-50/4	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	66.64	65.08	68.24
			Distill	72.01	70.31	72.82

CASIA-WebFace		Teacher Network		Student Network		
Top1 acc(%)	ResNet-50	Network	Initialization	Targets		
	74.16			ResNet-50/8	$\alpha = 0, \beta = 0$	$\alpha = 0, \beta = 1$
			Pretrain	51.86	51.43	53.45
			Distill	57.66	56.87	57.78

282 6 Conclusion

283 In this paper, we take face recognition as a breaking point, and propose the knowledge distillation
 284 on two non-classification tasks, including face alignment and verification. We extend the previ-
 285 ous distillation framework by transferring the distilled knowledge from face classification to face
 286 alignment and verification. By selecting appropriate initializations and targets, the distillation on
 287 non-classification tasks can be much easier. Besides, we also give some guidelines for target selection
 288 on non-classification tasks, and we hope these guidelines can be helpful for more tasks. Experiments
 289 on the CASIA-WebFace, CelebA and large-scale MS-Celeb-1M datasets have demonstrated the
 290 effectiveness of the proposed method, which gives the student networks that can be competitive or
 291 exceed the teacher network under appropriate compression rates. In addition, we use a common
 292 initialization trick to further improve the distillation performance of classification, and this can boost
 293 the distillation on non-classification tasks. Experiments on CASIA-WebFace have demonstrated the
 294 effectiveness of this simple trick.

295 **References**

- 296 [1] Lei Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *NIPS*, 2014.
- 297 [2] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise
298 training of deep networks. In *NIPS*, 2006.
- 299 [3] Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *KDD*,
300 2006.
- 301 [4] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An analysis of deep neural network
302 models for practical applications. In *arXiv:1605.07678*, 2017.
- 303 [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille.
304 Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*,
305 2015.
- 306 [6] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized
307 neural networks: Training deep neural networks with weights and activations constrained to +1
308 or -1. In *arXiv:1602.02830*, 2016.
- 309 [7] Emily Denton, Wojciech Zaremba, Yann LeCun, and Yann LeCun. Exploiting linear structure
310 within convolutional networks for efficient evaluation. In *NIPS*, 2014.
- 311 [8] Yunchao Gong, Liu Liu, Ming Yang, and Lubomir D. Bourdev. Compressing deep convolutional
312 networks using vector quantization. In *ICLR*, 2015.
- 313 [9] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset
314 and benchmark for large-scale face recognition. In *ECCV*, 2016.
- 315 [10] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural
316 networks with pruning, trained quantization and huffman coding. In *ICLR*, 2016.
- 317 [11] Song Han, Jeff Pool, Sharan Narang, Huizi Mao, Enhao Gong, Shijian Tang, Erich Elsen, Peter
318 Vajda, Manohar Paluri, John Tran, Bryan Catanzaro, and William J. Dally. Dsd: Dense-sparse-
319 dense training for deep neural networks. In *ICLR*, 2017.
- 320 [12] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections
321 for efficient neural networks. In *NIPS*, 2015.
- 322 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
323 recognition. In *CVPR*, 2016.
- 324 [14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In
325 *NIPS Workshop*, 2014.
- 326 [15] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep
327 belief nets. 2006.
- 328 [16] Forrest N. Iandola, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer.
329 Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5mb model size. In
330 *arXiv:1602.07360*, 2016.
- 331 [17] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training
332 by reducing internal covariate shift. In *ICML*, 2015.
- 333 [18] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
334 In *Technical report, University of Toronto*, 2009.
- 335 [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
336 convolutional neural networks. In *NIPS*, 2012.
- 337 [20] Erik Learned-Miller, Gary B. Huang, Aruni RoyChowdhury, and Gang Hua. Labeled faces in
338 the wild: A survey. In *Advances in Face Detection and Facial Image Analysis*, 2016.

- 339 [21] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
340 wild. In *ICCV*, 2015.
- 341 [22] Ping Luo, Zhenyao Zhu, Ziwei Liu, Xiaogang Wang, and Xiaoou Tang. Face model compression
342 by distilling knowledge from neurons. In *AAAI*, 2016.
- 343 [23] Alexander Novikov, Dmitry Podoprikin, Anton Osokin, and Dmitry Vetrov. Tensorizing neural
344 networks. In *NIPS*, 2015.
- 345 [24] George Papamakarios. Distilling model knowledge. In *arXiv:1510.02437*, 2015.
- 346 [25] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In *BMVC*,
347 2015.
- 348 [26] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet
349 classification using binary convolutional neural networks. In *ECCV*, 2016.
- 350 [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time
351 object detection with region proposal networks. In *NIPS*, 2015.
- 352 [28] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
353 Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.
- 354 [29] Bharat Bhuvan Sau and Vineeth N. Balasubramanian. Deep model compression: Distilling
355 knowledge from noisy teachers. In *arXiv:1610.09650*, 2017.
- 356 [30] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for
357 face recognition and clustering. In *CVPR*, 2015.
- 358 [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
359 image recognition. In *ICLR*, 2015.
- 360 [32] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov,
361 Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions.
362 In *CVPR*, 2015.
- 363 [33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Re-
364 thinking the inception architecture for computer vision. In *CVPR*, 2016.
- 365 [34] Cheng Tai, Tong Xiao, Yi Zhang, Xiaogang Wang, and Weinan E. Convolutional neural
366 networks with low-rank regularization. In *CoPR*, 2015.
- 367 [35] Gregor Urban, Krzysztof J. Geras, Samira Ebrahimi Kahou, Ozlem Aslan, Shengjie Wang,
368 Rich Caruana, Abdelrahman Mohamed, Matthai Philipose, and Matt Richardson. Do deep
369 convolutional nets really need to be deep and convolutional? In *arXiv:1603.05691*, 2017.
- 370 [36] Chong Wang, Xue Zhang, and Xipeng Lan. How to train triplet networks with 100k identities?
371 In *arXiv:1709.02940*, 2017.
- 372 [37] Yue Wu, Tal Hassner, KangGeon Kim, Gerard Medioni, and Prem Natarajan. Facial landmark
373 detection with tweaked convolutional neural networks. In *arXiv:1511.04031*, 2015.
- 374 [38] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z. Li. Learning face representation from scratch.
375 In *arXiv:1506.02640*, 2016.
- 376 [39] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in
377 deep neural networks? In *NIPS*, 2014.
- 378 [40] Xiangyu Zhang, Jianhua Zou, Kaiming He, and Jian Sun. Accelerating very deep convolutional
379 networks for classification and detection. *TPAMI*, 2016.
- 380 [41] Fang Zhao, Yongzhen Huang, Liang Wang, and Tieniu Tan. Deep semantic ranking based
381 hashing for multi-label image retrieval. In *CVPR*, 2015.