

Supplementary Materials

FaceLiVT: Face Recognition using Linear Vision Transformer with Structural Reparameterization

In this supplementary material, we present implementation details and additional quality comparisons to further demonstrate the effectiveness and quality of our proposed method in comparison to other existing approaches. By including these extended evaluations, we aim to provide a more comprehensive analysis and highlight the advantages of our method over alternative techniques.

A. Implementation Details of FaceLiVT

The advantages of FaceLiVT are using RepMixer, which can be reparameterized while inferring the model to speed up. Below is the implementation of RepMixer in PyTorch-like code:

```

1 class RepConv(nn.Module):
2     def __init__(self, inc, ouc, ks=1, stride=1,
3                 pad=0, groups=1):
4         super().__init__()
5         self.conv = nn.Conv2d(inc, ouc, ks, stride,
6                               pad, groups=groups)
7         self.repconv = nn.Conv2d(inc, ouc, ks//2,
8                                   stride, pad//2, groups=groups)
9         self.bn = nn.BatchNorm2d(ouc)
10
11     def forward(self, x):
12         xr = self.conv(x) + self.repconv(x)
13         return self.bn(xr)
14
15     def forward_deploy(self, x):
16         return self.conv(x)
17
18     @torch.no_grad()
19     def reparam(self):
20         conv = self.conv
21         repconv=self.repconv; self.__delattr__('repconv')
22         kw, kh = (conv.weight.shape[2]-
23                 repconv.weight.shape[2])//2,
24                 (conv.weight.shape[3]-
25                 repconv.weight.shape[3])//2
26         repconv_w = nn.functional.pad(repconv.weight,
27                                     [kh, kh, kw, kw])
28         repconv_b = repconv.bias
29
30         final_conv_w = conv.weight + repconv_w
31         final_conv_b = conv.bias + repconv_b
32
33         conv.weight.data.copy_(final_conv_w)
34         conv.bias.data.copy_(final_conv_b)
35
36         bn = self.bn
37         w = bn.weight / (bn.running_var + bn.eps)**0.5
38         w = conv.weight * w[:, None, None, None]
39         b = bn.bias + (conv.bias - bn.running_mean) *
40             bn.weight / (bn.running_var + bn.eps)**0.5
41         self.__delattr__('bn')
42
43         conv.weight.data.copy_(w)
44         conv.bias.data.copy_(b)
45
46         self.forward = self.forward_deploy
47         self.conv = conv
48         return self

```

The other advantages of FaceLiVT are using Multi-Head Linear Attention (MHLA), which can reduce the complexity of conventional Multi-Head Self Attention from quadratic $\mathcal{O}(N^2)$ into linear $\mathcal{O}(Nr)$. Below is the implementation of MHLA in PyTorch-like code:

```

1 class MultiHeadLinearAttention(torch.nn.Module):
2     def __init__(self, dim, resolution,
3                 ratio=4, act_layer=nn.ReLU):
4         super().__init__()
5         self.n_head = 16
6         self.res=resolution**2
7         self.dim=dim
8         lin1=[]
9         lin2=[]
10        self.norm=nn.GroupNorm(1, dim)
11        self.act_layer = act_layer()
12        for i in range(self.n_head):
13            lin1.append(nn.Linear(self.res, self.res*ratio))
14            lin2.append(nn.Linear(self.res*ratio, self.res))
15        self.lin1 = torch.nn.ModuleList(lin1)
16        self.lin2 = torch.nn.ModuleList(lin2)
17
18    def forward(self, x):
19        B,C,H,W= x.shape
20        x = self.norm(x).reshape(-1, self.dim, self.res)
21        x = x.chunk(self.n_head, dim=1)
22        x_out = []
23        for i in range(self.n_head):
24            feat = self.lin1[i](x[i])
25            feat = self.act_layer(feat)
26            feat = self.lin2[i](feat)
27            x_out.append(feat)
28        x = torch.cat(x_out, dim=-1)
29        x = x.reshape(B,C,H,W)
30        return x

```

We employ a similar training methodology as [1]–[3] on Glint360K [4]. The AdamW optimizer [5] and a polynomial decay learning rate schedule are utilized. The starting learning rate is 6×10^{-3} with a minimum of 1×10^{-5} . The total batch size amounts to 384 using $3 \times$ NVIDIA RTX-A6000 GPUs, and the weight decay is 1×10^{-4} . The model is trained for 20 epochs with a pre-processing resolution of (112×112) .

TABLE I
FACELiVT TRAINING HYPERPARAMETERS ON GLINT360K DATASETS

Hyperparameters	Config
optimizer	AdamW
learning rate	0.006
batch size	384
LR schedule	polynomial
warmup epochs	0
training epochs	20
weight decay	0.0001
embedding size	512

B. Detailed Result On Benchmark Test

Table II shows the extension detail of testing in the IJB-B and IJB-C datasets with TAR at different FAR values from $1e^{-5}$ to $1e^{-3}$. Our proposed model is slightly lower than EdgeFace, since it was trained using a larger dataset, WebFace240M, and 80 epochs longer than ours. Due to the limited computation hardware, we trained using Glint360K. With the similar datasets, our proposed FaceLiVT can outperform the CNN-based face recognition model.

Table III presents the inference speed of our proposed FaceLiVT. Although EdgeFace, a Hybrid-ViT with reduced

TABLE II
COMPARISON OF FACELiVT VARIANT WITH STATE-OF-THE-ART ON
IJB-B AND IJB-C BENCHMARK DATASET.

Model	IJB-B			IJB-C		
	$1e^{-5}$	$1e^{-4}$	$1e^{-3}$	$1e^{-5}$	$1e^{-4}$	$1e^{-3}$
MobileFaceNetV1 [2]	85.7	92.0	94.7	90.4	93.9	95.9
MobileFaceNet [2]	87.9	92.8	95.6	92.2	94.7	96.6
ShuffleFaceNet-1.5 [2]	86.5	92.3	95.2	91.3	94.3	96.3
EdgeFace-S(0.5) [1]	-	93.6	-	-	95.6	-
ResNet50-ArcFace [6]	80.5	89.9	94.5	86.1	92.1	96.0
FaceLiVT-M	88.4	93.4	95.6	91.3	95.0	96.7
FaceLiVT-M-(LA)	83.7	92.5	95.1	89.25	94.1	96.3
ShuffleFaceNet-0.5 [2]	-	-	-	-	-	-
EdgeFace-XS(0.6) [1]	-	92.7	-	-	94.8	-
FaceLiVT-S	68.8	89.1	95.0	66.7	89.7	96.0
FaceLiVT-S-(LA)	58.4	83.4	94.7	56.3	82.5	94.63

parameters and FLOPs, presents a bottleneck in the mobile latency inference test. In comparison to EdgeFace-XS, our proposed FaceLiVT-M(LA) is 8.6 times faster, despite having higher FLOPs, due to effective reparameterization.

TABLE III
COMPARISON OF FACELiVT VARIANT WITH STATE-OF-THE-ART ON
IPHONE 15 PRO MOBILE LATENCY TEST.

Model	Type	Param (M)	FLOP (M)	Latency (ms)
MobileFaceNetV1 [2]	CNN	3.4	1100	0.81
MobileFaceNet [2]	CNN	2.0	933	0.77
ShuffleFaceNet-1.5 \times [2]	CNN	2.6	577	0.69
EdgeFace-S(0.5) [1]	Hybrid ViT	3.6	306	10.21
FaceLiVT-M	Hybrid ViT	14.3	569	1.11
FaceLiVT-M-(LA)	Hybrid ViT	9.75	386	0.67
ShuffleFaceNet-0.5 [2]	CNN	1.4	66.9	0.45
EdgeFace-XS(0.6) [1]	Hybrid ViT	1.77	154	5.82
FaceLiVT-S	Hybrid ViT	5.89	237	0.61
FaceLiVT-S-(LA)	Hybrid ViT	5.05	160	0.47

REFERENCES

- [1] A. George, C. Ecabert, H. O. Shahreza, K. Kotwal, and S. Marcel, "Edgeface: Efficient face recognition model for edge devices," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2024.
- [2] Y. Martinez-Diaz, M. Nicolas-Diaz, H. Mendez-Vazquez, L. S. Luevano, L. Chang, M. Gonzalez-Mendoza, and L. E. Sucar, "Benchmarking lightweight face architectures on specific face recognition scenarios," *Artificial Intelligence Review*, vol. 54, no. 8, pp. 6201–6244, 2021.
- [3] Y. Martindez-Diaz, L. S. Luevano, H. Mendez-Vazquez, M. Nicolas-Diaz, L. Chang, and M. Gonzalez-Mendoza, "Shufflefacenet: A lightweight face architecture for efficient and highly-accurate face recognition," in *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pp. 0–0, 2019.
- [4] X. An, X. Zhu, Y. Gao, Y. Xiao, Y. Zhao, Z. Feng, L. Wu, B. Qin, M. Zhang, D. Zhang, *et al.*, "Partial fc: Training 10 million identities on a single machine," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1445–1449, 2021.
- [5] I. Loshchilov, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [6] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4690–4699, 2019.