

5 Supplementary

Algorithm 1 PyTorch Style Pseudocode for Progressive Knowledge Distillation

```

1: Initialize models:  $S, \{T_1, T_2, T_3, T_4\}$ 
2: Define:  $N$  maxiterations,  $C$  interval,  $patience$ 
3: DataLoader  $D_{train}$ ,  $\lambda$ ,  $iter \leftarrow 0$ ,  $T \leftarrow T_1$ ,  $pat \leftarrow 0$ ,  $\mathcal{L}_{KD}^{prev} \leftarrow \infty$ 
4: while  $iter < N$  and  $perf < \tau$  do
5:    $iter \leftarrow iter + 1$ 
6:   for batch  $X$ , targets  $Y$  in  $D_{train}$  do
7:      $S.train, T.eval$ 
8:      $P_T \leftarrow T(X_{high\ res}), P_S \leftarrow S(X_{low\ res})$ 
9:      $\mathcal{L}_{GT} \leftarrow Loss(P_S, Y)$ 
10:     $\mathcal{L}_{KD} \leftarrow Loss(S_{bottleneck}, T_{bottleneck}) + Loss(P_S, softmax(P_T))$ 
11:     $\mathcal{L}_{Total} \leftarrow \mathcal{L}_{GT} + \lambda \cdot \mathcal{L}_{KD}$ 
12:    Update  $S$  using  $\mathcal{L}_{Total}$ 
13:    if  $iter \bmod C == 0$  or  $(\mathcal{L}_{KD} \geq \mathcal{L}_{KD}^{prev}$  and  $pat \geq patience)$  then
14:       $T \leftarrow$  Next teacher in  $\{T_1, T_2, T_3, T_4\}$ 
15:       $pat \leftarrow 0$ 
16:    else
17:       $pat \leftarrow pat + 1$ 
18:    end if
19:     $\mathcal{L}_{KD}^{prev} \leftarrow \mathcal{L}_{KD}$ 
20:  end for
21: end while

```

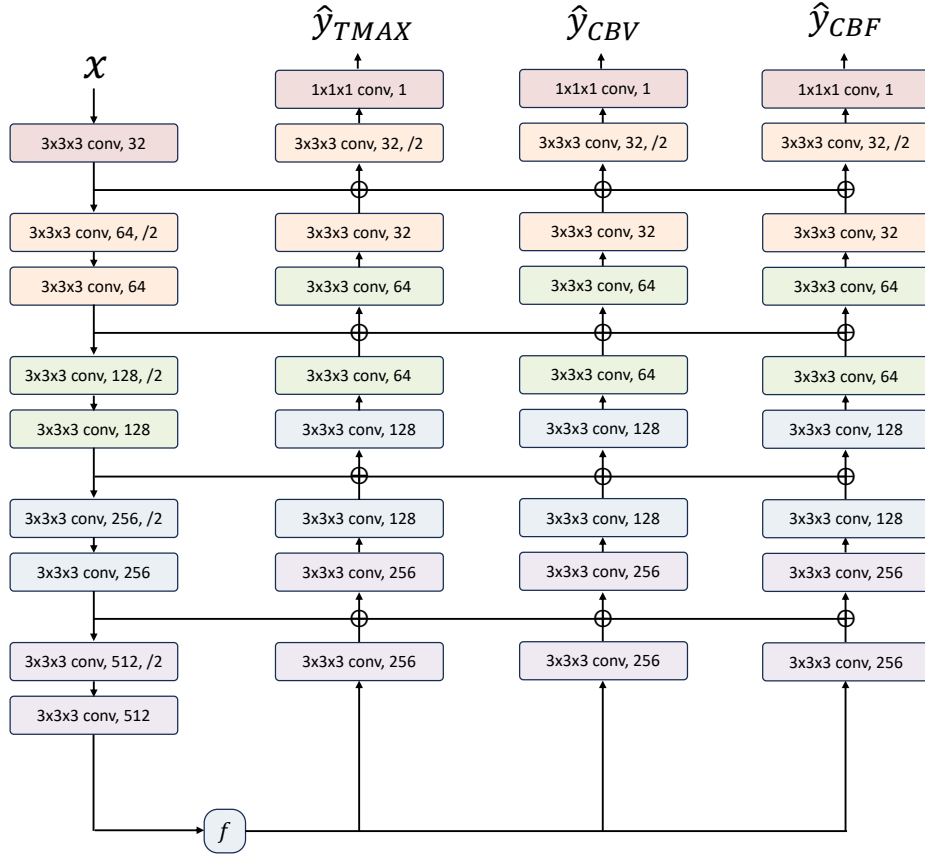


Fig. 3. Detailed architecture of the components in our framework. Intermediate features and bottleneck feature f of the encoder are used as an additional input to each TMAX, CBV, and CBF task specific decoders.