

## 6 Appendix

### 6.1 Full experiment results

In addition to our main paper examples, we test our method across a diverse range of 44 pre-trained models. These models vary in scale, architecture, development year, and purpose. Our method consistently reduces computation significantly in all cases, with minimal or no loss in accuracy.

**Pool of models** Models in our experiments include:

1. ConvNeXtv2 [66]: *atto, femto, pico, nano, tiny, base, large, huge*
2. DaViT [14]: *tiny, small, base*
3. MaxViT [59]: *tiny, small, base, large*
4. SwinT [41]: *tiny, small, base, large*
5. ResNet [25]: *18, 34, 50, 101, 152*
6. EfficientNet [57]: *b0, b1, b2, b3, b4*
7. EfficientViT [2]: *b0, b1, b2, b3*
8. EfficientViT [40]: *m1, m2, m3, m4, m5*
9. EfficientFormerv2 [35]: *s0, s1, s2, l*
10. MobileNetv3 [27]: *small100, large100*

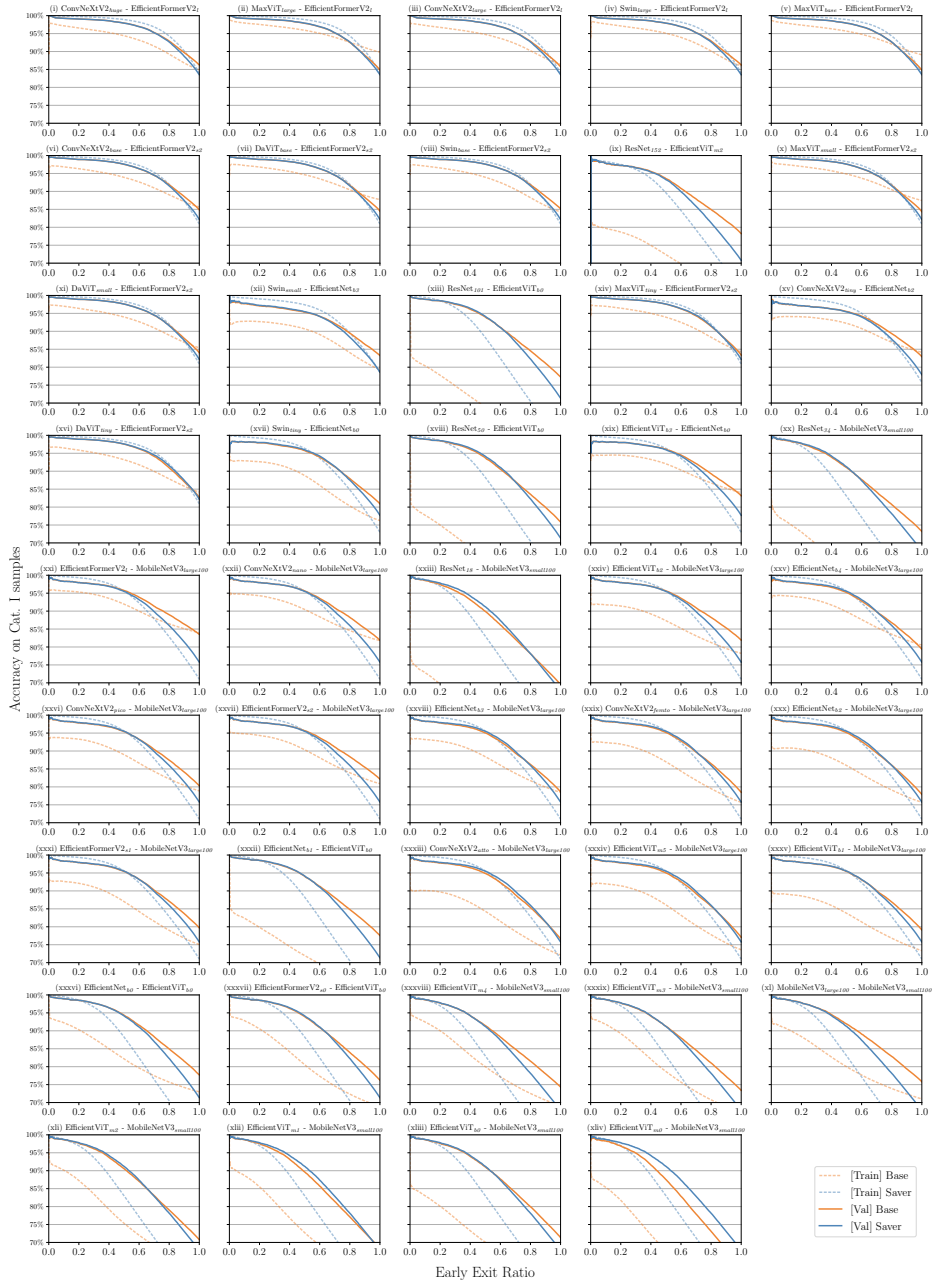
All these models are trained in their original ways with the input size of  $224 \times 224$  for the ImageNet-1k classification task. Some models are pre-trained with larger datasets [66]. The model weights of ResNet is provided by torchvision [44], all other models are provided by timm [65] and their complexities are measured by the profiling tool of DeepSpeed [49]. Each model can employ a smaller one as its saver, as described in Sec. 3.3.

All models were trained using their original procedures with an input size of  $224 \times 224$  for the ImageNet-1k classification task. Some models are pre-trained on larger datasets [66]. ResNet weights are available from torchvision [44], while weights for other models are accessed via timm [65]. We assessed their complexities using DeepSpeed’s profiling tool [49]. According to the methodology described in Sec. 3.3, each model could utilize a smaller one as its saver.

**Compression effect** Tab. 5 complements Tab. 1, providing extensive experimental results that reinforce our primary findings discussed in Sec. 4.1. Typically, larger models exhibit more substantial compression potential and can maintain their performance even after significant reductions. In contrast, compressing smaller models presents challenges, as their corresponding saver models are relatively larger in proportion. The smallest models, EfficientViT<sub>m0</sub> and MobileNetv3<sub>small100</sub>, are not used as base models due to the absence of smaller eligible savers. Moreover, for effective zero-loss compression, saver models need to closely approximate the accuracy of their base models; otherwise, a trade-off between performance and compression is inevitable.

Base Model Name	GFLOPs.	Acc.	Saver Model			Required MACs(%) when Acc. Drop $\leq$					Max Performance	
			Name	GFLOPs.	Acc.	1%	0.5%	0.01%	0%	Avg. GMACs	Acc.	
ConvNeXtV2 <sub>huge</sub> [66]	114.92	86.25%	EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	13.21	19.41	31.17	67.34	70.70	86.25%	
						(-88.5%)	(-83.1%)	(-72.9%)	(-41.4%)	(-38.5%)	(-0.00%)	
MaxViT <sub>large</sub> [59]	42.80	84.82%	EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	3.46	5.48	7.81	9.04	11.24	84.92%	
						(-91.9%)	(-87.2%)	(-81.8%)	(-78.9%)	(-73.7%)	(+0.09%)	
ConvNeXtV2 <sub>large</sub> [66]	34.36	85.76%	EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	5.16	7.15	11.00	15.15	17.49	85.78%	
						(-85.0%)	(-79.2%)	(-68.0%)	(-55.9%)	(-49.1%)	(+0.02%)	
SwiT <sub>large</sub> [41]	34.02	86.24%	EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	5.56	7.22	10.95	14.46	18.70	86.26%	
						(-83.7%)	(-78.8%)	(-67.8%)	(-57.5%)	(-45.0%)	(+0.02%)	
MaxViT <sub>base</sub> [59]	23.39	84.80%	EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	2.95	4.13	5.33	6.74	7.28	84.86%	
						(-87.4%)	(-82.3%)	(-77.2%)	(-71.2%)	(-68.9%)	(+0.06%)	
ConvNeXtV2 <sub>base</sub> [66]	15.35	84.87%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	2.74	3.57	4.85	5.44	6.35	84.92%	
						(-82.1%)	(-76.8%)	(-68.4%)	(-64.6%)	(-58.6%)	(+0.05%)	
DaViT <sub>base</sub> [14]	15.22	84.48%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	2.40	3.16	4.40	5.19	6.50	84.53%	
						(-84.3%)	(-79.2%)	(-71.1%)	(-65.9%)	(-57.3%)	(+0.04%)	
SwiT <sub>base</sub> [41]	15.13	85.14%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	2.81	3.70	4.79	9.06	9.40	85.15%	
						(-81.4%)	(-75.6%)	(-68.3%)	(-40.1%)	(-37.9%)	(+0.00%)	
ResNet <sub>102</sub> [25]	11.51	78.24%	EfficientViT <sub>m2</sub> [40]	0.20	70.78%	4.17	4.91	6.30	8.79	9.92	78.27%	
						(-63.8%)	(-57.3%)	(-45.3%)	(-23.7%)	(-13.8%)	(-0.03%)	
MaxViT <sub>small</sub> [59]	11.31	84.33%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	2.01	2.51	3.53	3.86	4.29	84.42%	
						(-82.2%)	(-77.8%)	(-68.8%)	(-65.9%)	(-62.1%)	(-0.08%)	
DaViT <sub>small</sub> [14]	8.58	84.00%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	1.65	2.06	2.92	3.25	3.54	84.06%	
						(-80.8%)	(-76.0%)	(-66.0%)	(-62.1%)	(-58.7%)	(+0.05%)	
SwiT <sub>small</sub> [41]	8.51	83.25%	EfficientNet <sub>s3</sub> [57]	0.96	78.63%	2.58	3.11	4.09	4.41	5.67	83.34%	
						(-69.6%)	(-63.5%)	(-51.9%)	(-48.2%)	(-33.4%)	(+0.09%)	
ResNet <sub>101</sub> [25]	7.80	77.26%	EfficientViT <sub>h0</sub> [2]	0.10	71.35%	2.10	2.53	2.93	3.46	4.43	77.38%	
						(-73.1%)	(-67.5%)	(-62.5%)	(-55.7%)	(-43.3%)	(+0.12%)	
MaxViT <sub>tiny</sub> [59]	5.37	83.36%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	1.29	1.52	1.82	1.96	2.57	83.58%	
						(-76.0%)	(-71.7%)	(-66.0%)	(-63.5%)	(-52.2%)	(+0.22%)	
ConvNeXtV2 <sub>tiny</sub> [66]	4.46	82.94%	EfficientNet <sub>s3</sub> [57]	0.66	77.89%	1.65	1.98	2.54	3.50	4.54	82.96%	
						(-62.9%)	(-55.6%)	(-42.9%)	(-21.4%)	(+2.0%)	(+0.01%)	
DaViT <sub>tiny</sub> [14]	4.41	82.71%	EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	-	1.23	1.37	1.42	2.06	83.23%	
						-	(-72.1%)	(-68.9%)	(-67.9%)	(-53.3%)	(+0.51%)	
SwiT <sub>tiny</sub> [41]	4.35	80.89%	EfficientNet <sub>s0</sub> [57]	0.39	77.70%	0.91	1.15	1.42	1.54	2.07	81.09%	
						(-79.0%)	(-73.5%)	(-67.4%)	(-64.5%)	(-52.3%)	(+0.19%)	
ResNet <sub>50</sub> [25]	4.09	75.85%	EfficientViT <sub>h0</sub> [2]	0.10	71.35%	0.93	1.11	1.31	1.36	1.87	76.15%	
						(-77.2%)	(-72.9%)	(-67.9%)	(-66.7%)	(-54.2%)	(+0.30%)	
EfficientViT <sub>h3</sub> [2]	3.88	83.14%	EfficientNet <sub>h0</sub> [57]	0.39	77.70%	1.33	1.59	2.08	3.17	3.89	83.16%	
						(-65.6%)	(-59.0%)	(-46.4%)	(-18.4%)	(+0.3%)	(+0.01%)	
ResNet <sub>34</sub> [25]	3.66	73.2%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	1.11	1.35	1.57	1.63	2.01	73.36%	
						(-69.7%)	(-63.0%)	(-57.3%)	(-55.6%)	(-45.1%)	(+0.16%)	
EfficientFormerV2 <sub>t</sub> [35]	2.51	83.53%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.99	1.17	1.49	2.51	2.51	83.53%	
						(-60.6%)	(-53.4%)	(-40.6%)	(+0.1%)	(+0.1%)	(+0.00%)	
ConvNeXtV2 <sub>nano</sub> [66]	2.45	81.87%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.88	1.04	1.34	1.76	2.45	81.87%	
						(-63.9%)	(-57.5%)	(-45.2%)	(-28.2%)	(+0.3%)	(+0.00%)	
ResNet <sub>18</sub> [25]	1.81	69.53%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.19	0.28	0.37	0.41	0.84	70.31%	
						(-89.4%)	(-84.3%)	(-79.5%)	(-77.2%)	(-54.0%)	(+0.78%)	
EfficientViT <sub>b2</sub> [2]	1.55	81.91%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.64	0.74	0.95	1.65	1.65	81.91%	
						(-58.9%)	(-52.5%)	(-38.6%)	(+6.2%)	(+6.2%)	(+0.00%)	
EfficientNet <sub>h4</sub> [57]	1.50	79.40%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.46	0.54	0.64	0.66	0.78	79.66%	
						(-69.4%)	(-63.8%)	(-57.7%)	(-56.1%)	(-48.4%)	(+0.26%)	
ConvNeXtV2 <sub>pico</sub> [66]	1.37	80.31%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.52	0.61	0.71	0.81	1.05	80.38%	
						(-61.9%)	(-55.4%)	(-48.0%)	(-40.6%)	(-23.1%)	(+0.07%)	
EfficientFormerV2 <sub>s2</sub> [35]	1.22	82.15%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.57	0.64	0.82	1.02	1.05	82.17%	
						(-53.2%)	(-47.2%)	(-33.0%)	(-16.2%)	(-13.8%)	(+0.01%)	
EfficientNet <sub>h3</sub> [57]	0.96	78.63%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.33	0.38	0.43	0.45	0.55	79.01%	
						(-66.0%)	(-61.0%)	(-55.4%)	(-53.2%)	(-42.4%)	(+0.38%)	
ConvNeXtV2 <sub>gen2o</sub> [66]	0.78	78.49%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.31	0.36	0.42	0.43	0.51	78.76%	
						(-60.2%)	(-53.4%)	(-46.7%)	(-44.6%)	(-34.9%)	(+0.27%)	
EfficientNet <sub>h2</sub> [57]	0.66	77.89%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.27	0.30	0.34	0.36	0.46	78.29%	
						(-59.4%)	(-53.9%)	(-47.8%)	(-45.5%)	(-30.9%)	(+0.39%)	
EfficientFormerV2 <sub>s1</sub> [35]	0.63	79.69%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.34	0.38	0.43	0.50	0.58	79.76%	
						(-46.7%)	(-39.9%)	(-31.4%)	(-21.4%)	(-8.6%)	(+0.06%)	
EfficientNet <sub>h1</sub> [57]	0.57	77.57%	EfficientViT <sub>h0</sub> [2]	0.10	71.35%	0.26	0.28	0.33	0.35	0.38	77.64%	
						(-54.9%)	(-50.2%)	(-42.4%)	(-38.8%)	(-33.3%)	(+0.06%)	
ConvNeXtV2 <sub>attno</sub> [66]	0.55	76.64%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	-	0.23	0.26	0.26	0.39	77.33%	
						-	(-57.6%)	(-52.8%)	(-51.7%)	(-29.6%)	(+0.68%)	
EfficientViT <sub>m5</sub> [40]	0.52	77.08%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.23	0.25	0.28	0.29	0.38	77.55%	
						(-55.7%)	(-51.2%)	(-44.9%)	(-43.0%)	(-27.1%)	(+0.47%)	
EfficientViT <sub>b1</sub> [2]	0.51	79.12%	MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	0.30	0.34	0.37	0.40	0.45	79.20%	
						(-40.8%)	(-33.4%)	(-26.8%)	(-22.6%)	(-11.5%)	(+0.08%)	
EfficientNet <sub>h0</sub> [57]	0.39	77.70%	EfficientViT <sub>h0</sub> [2]	0.10	71.35%	0.21	0.23	0.29	0.33	0.36	77.71%	
						(-45.6%)	(-39.5%)	(-25.2%)	(-14.6%)	(-7.2%)	(+0.01%)	
EfficientFormerV2 <sub>s0</sub> [35]	0.38	76.24%	EfficientViT <sub>h0</sub> [2]	0.10	71.35%	0.18	0.20	0.22	0.23	0.26	76.38%	
						(-51.0%)	(-46.0%)	(-41.3%)	(-38.8%)	(-30.6%)	(+0.13%)	
EfficientViT <sub>m4</sub> [40]	0.30	74.32%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.16	0.18	0.21	0.23	0.26	74.34%	
						(-47.3%)	(-40.3%)	(-29.2%)	(-21.6%)	(-12.9%)	(+0.02%)	
EfficientViT <sub>m3</sub> [40]	0.26	73.38%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.14	0.16	0.18	0.19	0.21	73.45%	
						(-46.2%)	(-39.7%)	(-30.8%)	(-26.6%)	(-18.3%)	(+0.06%)	
MobileNetV3 <sub>large100</sub> [27]	0.22	75.78%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.14	0.16	0.19	0.22	0.22	75.78%	
						(-33.6%)	(-26.5%)	(-12.2%)	(+2.8%)	(+3.4%)	(+0.00%)	
EfficientViT <sub>m2</sub> [40]	0.20	70.78%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.09	0.10	0.12	0.12	0.15	71.13%	
						(-54.5%)	(-48.3%)	(-40.9%)	(-38.8%)	(-22.6%)	(+0.35%)	
EfficientViT <sub>m1</sub> [40]	0.16	68.32%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	-	0.06	0.06	0.07	0.11	69.41%	
						-	(-64.1%)	(-60.7%)	(-59.6%)	(-31.9%)	(+1.09%)	
EfficientViT <sub>h0</sub> [2]	0.10	71.35%	MobileNetV3 <sub>small100</sub> [27]	0.06	67.65%	0.08	0.09	0.10	0.10	0.11	71.46%	
						(-16.5%)	(-8.7%)	(-0.6%)	(+3.7%)	(+11.5%)	(+0.10%)	

Table 5: Compression effect while TinySaver applied to all models in the pool



**Fig. 6: Accuracy vs early exit ratio on the ImageNet-1k training/validation set.** The intersection of solid lines denotes the ratio where the saver and base model have the equivalent accuracy on the validation set. The intersection of dashed lines are for the training set. Every plot is labelled as **Base - Saver**.

**Performance and early exit ratio** Fig. 6 expands on Fig. 3, illustrating how accuracy varies as more samples are processed by the saver model. In all test cases, the saver model initially excels on the training set. Cat. I samples are filtered based on saver’s confidence and the saver can perform exceptionally well with high-confidence samples. While base models struggle with these same samples, leading to poorer performance when Cat. I is small. Additionally, data augmentation complicates early performance for base models. However, as Cat. I expands to include more samples, the situation changes. Finally, if we set saver’s threshold as 0 and let all samples be Cat. I, the base will apparently outperform the saver.

On the validation set, the significant initial advantage of the saver model decreases, but it remains more accurate until a later intersection point. This consistent trend across all plots demonstrates the reliability of models when they are highly confident, which is the source of compression effect.

Nonetheless, the base model’s underperformance on the training set results in a relative overconfidence in the saver model. Although the validation set is not directly impacted, this discrepancy could adversely affect the accuracy of saver selection, particularly in multi-exit systems like ESN, which rely on exit policies derived from training data. Introducing more data augmentation on the training set may help reduce the discrepancy between training and validation performances.

## 6.2 Speed test

In our experiments, we use original pre-trained models, which are implemented differently by each author. Their heterogeneity makes them unevenly supported on different platforms. Therefore, it’s difficult to compare FLOPs reduction and actual speedups directly. However, for a given deployment scenario, we can replace FLOPs with measured latency while identifying the best saver model. The speed test on a PC with i9-11900F and one RTX 3090 is below. We ran inference on IN-1k val set, and recorded the model throughput. All TinySaver enabled models were configured to have no accuracy loss. Tab. 6 shows the results. Though models exhibit varying efficiency on different deployment scenarios, we can still observe that TinySaver can significantly improve the throughput of the base model.

## 6.3 Exit Sequence Network (ESN)

The design of the ESN comprises three components: (1) Embedding units that transform features at various scales into a format compatible with the subsequent sequential model. (2) A sequential feature bus designed to aggregate features from all preceding steps and transmit them to subsequent stages. (3) Head units, responsible for outputting results like regular attached exits.

To ensure efficiency, the overall scale of the ESN should be kept moderate, thereby minimizing the introduction of overhead. While over-simplified networks

Base Model	Saver			
	Original→Ours Avg. throughput img/s			
	bs=1,CPU	bs=1	bs=128	bs=512
ConvNextv2 <sub>h</sub>	Swin <sub>b</sub> 1.9→6.4	ConvNextv2 <sub>l</sub> 55.0→98.5	Swin <sub>b</sub> 96.9→307.8	OOM
Swin <sub>b</sub>	EfficientFormerV2 <sub>l</sub> 9.9→17.7	ConvNextv2 <sub>b</sub> 97.2→117.5	DaViT <sub>t</sub> 487.0→659.7	DaViT <sub>t</sub> 496.3→688.8
EfficientNet <sub>b2</sub>	ConvNextv2 <sub>a</sub> 51.9→81.4	ConvNextv2 <sub>a</sub> 150.9→229.4	EfficientViT <sub>m5</sub> 2172.1→2649.7	EfficientViT <sub>m5</sub> 2237.5→3344.4
EfficientFormerV2 <sub>s2</sub>	ConvNextv2 <sub>n</sub> 34.0→42.0	ConvNextv2 <sub>n</sub> 83.8→222.3	EfficientViT <sub>b2</sub> 562.1→1339.9	EfficientViT <sub>b2</sub> 570.8→1542.8

**Table 6:** Average throughput of models with TinySaver enabled on different deployment scenarios

might not be able to keep its functions. Careful tuning should be required for the best practice of ESN. We illustrate the detail of ESN in Fig. 7

**Feature embedding** In our approach, both saver and base models are considered well-trained and remain unchanged during ESN training. Their features, treated as a sequence of tokens, need embedding prior to joining the attention layer. In our ESN design, we employ small networks each containing a global average pooling layer and a linear layer, to embed intermediate backbone features.

**Sequential feature bus** The primary objective of this part is to ensure the delivery of the most informative feature for each head unit. The sequential execution is mandatory, aligning with the step-by-step nature of the dynamic inference process. For this purpose, we have chosen to implement a transformer decoder network recognized for its effectiveness across a broad of applications. Meanwhile, the first step, i.e. features from the saver model, are critical in our system. The transformer allows for direct referencing of these features without requiring intermediary states. This is another reason to select transformers. Contrasting with the commonly used Feature Pyramid Network (FPN) [37], our Exit Sequence Network (ESN) does not accept backward connections from deeper layers. This is required by EE to inference progressively in practice. Therefore, our implementation uses mask to maintain causality.

**Exiting heads** Exiting heads are responsible for generating the output. We also brings the saver model’s output, making  $\tilde{\mathbf{y}}_n = \tilde{\mathbf{y}}_S + \Delta\tilde{\mathbf{y}}_n$ . The weight of ESN heads are initialized as zero and only need to learn  $\Delta\tilde{\mathbf{y}}_n$ . Therefore, the intermediate outputs combine the information from the saver and base model and their performance is lower bounded by the saver model.

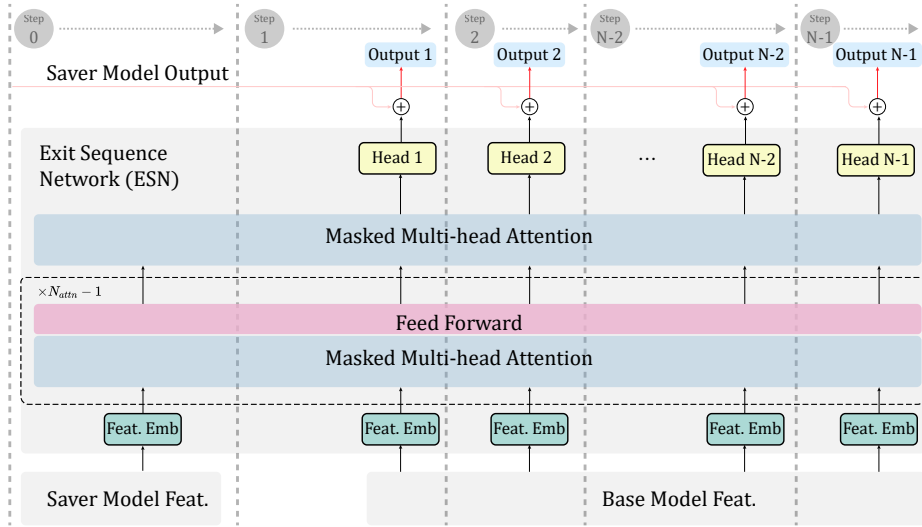


Fig. 7: Detailed ESN

### 6.4 Configuration and training details of ESN

Fig. 7 shows detailed structure of our prototype design of the ESN. We extract intermediate features from each residual block’s output, integrating them with saver features into the ESN of multi-head causal attention layers. Subsequently, single-layer classifiers at each exit generate predictions, all supervised by the loss function detailed in Sec. 6.4. Tab. 7 includes the hyper-parameter of training ESN. We run a basic grid search of suitable attention layer numbers and dims.

Fig. 7 illustrates the detailed structure of our ESN prototype. We extract intermediate features from the output of each backbone segment and combine them with features from the saver into the ESN, which includes multi-head causal attention layers. Then, single-layer classifiers at each exit are responsible for generating predictions, trained by the loss function detailed in Sec. 6.4. Additionally, Tab. 7 lists the hyperparameters used for training the ESN, where we conducted a basic grid search to identify optimal numbers of attention layers and their dimensions.

**Loss function** In traditional Early Exit (EE) models, the training usually employs the original loss function or knowledge distillation. However, in our approach, a significant portion of inputs is processed solely by the saver model, meaning intermediate exits do not influence these particular inputs. Additionally, lightweight networks have limited capacity for knowledge representation, suggesting the utility of reducing the burden on such networks from samples unlikely to reach intermediate exits. Consequently, we adjust the weight of samples in the loss function based on the saver model’s confidence levels. The loss

config	value
optimizer	AdamW
learning rate	$6.25 \times 10^{-4}$
learning rate schedule	cosine decay
warmup epochs	20
augmentation	AutoAugment [9]
mixup	0.8
cutmix	1.0
training epochs	10
num attn heads	8
num attn layers	{0, 1, 2, 3}
dim	{8, 32, 128, 512}

**Table 7:** Configurations for training ESN

function applied at the  $n^{th}$  exit is outlined in Eq. (11), reflecting this adaptation to reduce the load on intermediate exits while maintaining effective learning.

$$L_n = \sum_{m=0}^M 2(1 - \max_i(\tilde{\mathbf{y}}_S^{(m,i)})) \mathcal{L}(\tilde{\mathbf{y}}_n^{(m)}) \quad (11)$$

Where  $\max_i(\tilde{\mathbf{y}}_S^{(m,i)})$  is the confidence of the saver model.  $\mathcal{L}(\tilde{\mathbf{y}}_n^{(m)})$  is the regular loss function per sample and we use cross entropy in the experiment of this paper.