a) Pruner Training Processing **Block Importance** SA / FC Layers **Scores** $\times L$ $(T \times N)$ D_3^{T-1} Head $\mathcal{L}(X_T; \Theta, W)$ Queries D_3^{T-2} \overline{D} SA / FC FC Optimization Prediction Loss Prune (SA **Forward** 0.26 U_3^{T-2} **Backward** Once trained b) Pruning Post Processing Binary search initialization **Downsample Blocks** 1. Obtain current threshold c_t **Upsample Blocks** 2. Update Sates of Blocks **Pruned Blocks** 3. Calculate pruned ratio P* 4. Judge abs(P*-P) < 0.0125**X** Remove Dependences **Expected** 5. Update c_t to step 1 or stop **Attached New Paths Pruning Ratio▶** Ignored Paths Binary search end