











T/1024 time steps, 2048 freq.

$$ZDecoder_1(C_{in} = 48, C_{out} = 4 \cdot 2 \cdot 2)$$

T/1024 time steps, 512 freq.

$$ZDecoder_2(C_{in} = 96, C_{out} = 48)$$

T/1024 time steps, 128 freq.

$$ZDecoder_3(C_{in} = 192, C_{out} = 96)$$

T/1024 time steps, 32 freq.

$$ZDecoder_4(C_{in} = 384, C_{out} = 192)$$

T/1024 time steps, 8 freq.

T time steps.

$$TDecoder_1(C_{in} = 48, C_{out} = 4 \cdot 2)$$

T/4 time steps

$$TDecoder_2(C_{in} = 96, C_{out} = 48)$$

T/16 time steps

$$TDecoder_3(C_{in} = 192, C_{out} = 96)$$

T/64 time steps

$$TDecoder_4(C_{in} = 384, C_{out} = 192)$$

T/256 time steps

Cross-Domain Transformer Encoder

T/1024 time steps, 8 freq.

 $ZEncoder_4(C_{in} = 192, C_{out} = 384)$

T/1024 time steps, 32 freq.

 $ZEncoder_3(C_{in} = 96, C_{out} = 192)$

T/1024 time steps, 128 freq.

 $ZEncoder_2(C_{in} = 48, C_{out} = 96)$

 $ZEncoder_1(C_{in} = 2 \cdot 2, C_{out} = 48)$

T/1024 time steps, 512 freq.

T/1024 time steps 2048 freq.

T/256 time steps

 $TEncoder_4(C_{in} = 192, C_{out} = 384)$

T/64 time steps

 $TEncoder_3(C_{in} = 96, C_{out} = 192)$

T/16 time steps

 $\mathrm{TEncoder}_2(C_{in}=48,C_{out}=96)$

T/4 time steps

 $\mathrm{TEncoder}_1(C_{in}=2,C_{out}=48)$

T time steps

4 2

STFT

Monthe