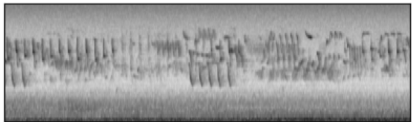


Audio Spectrogram

$\mathbf{x}_{0,\dots,T}$



\mathcal{E}_{CNN}

Encoder



Framed Feature Maps

\mathcal{E}_{CNN}

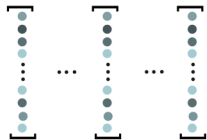


Aggregated in
Frequency

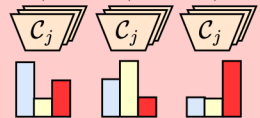
$\mathcal{E}_{\text{MLP}}^z$

Latent Representations

$q(\mathbf{z}_{0,\dots,T} | \mathbf{x}_{0,\dots,T})$



Downstream Classification



Intra-Frame Shift

$\hat{\delta}_{0\dots T}$

