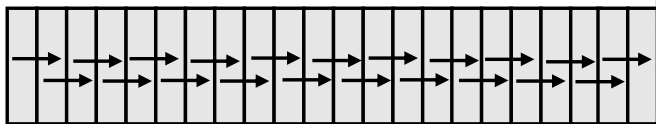
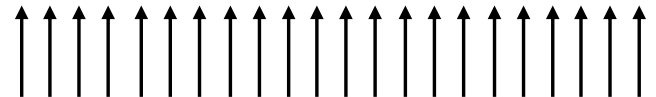


Context model  
(e.g., RNN)



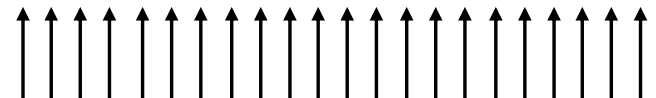
$\mathbf{z}_t$



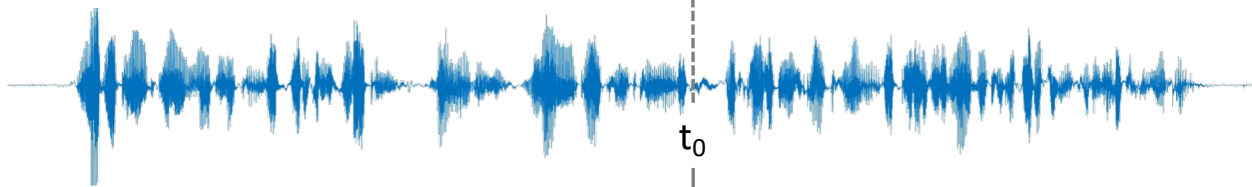
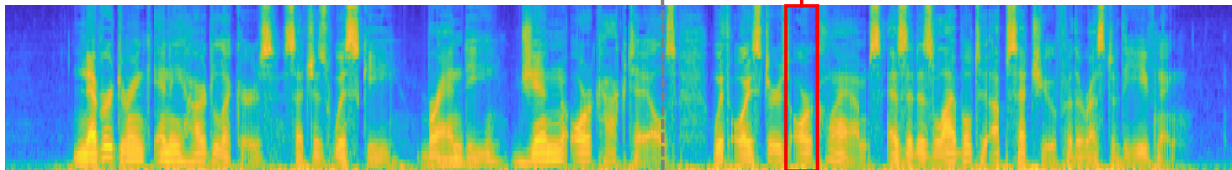
Frame-level  
encoder  
(e.g., MLP)



$\mathbf{y}_t$



log-Mel  
feature  
extraction



time

← past context | future input →

$\mathbf{c}_{t0}$



linear  
projection

$\mathbf{y}_{t+k}^*$

prediction  
error

$\mathbf{y}_{t+k}$

$$\mathcal{L}(\mathbf{y}_{t+k}^*, \mathbf{y}_{t+k}) = |\mathbf{y}_{t+k}^* - \mathbf{y}_{t+k}|$$