# SepFormer: High-Level



X → Encoder → h → Masking Net → $m_1$, $m_2$ → ⊗ → Decoder → $\hat{s}_1$, $\hat{s}_2$

## Encoder:

$$h = ReLU(\,conv1d(x)\,)$$

$$\begin{cases} X \in \mathbb{R}^T \\ h \in \mathbb{R}^{F \times T'} \end{cases}$$

## Masking Network:

h → Masking Net → $\{m_1, ..., m_N\}$ ; N masks, one for each speaker

h → Layer Norm → Linear→F → h' → Chunking → h'' → Sep Former → h''' → PreLN+Linear → h'''' → Overlay Add → h''''' → PReLU+FFW → $m_1$, $m_2$

## Chunking:

"create overlap chunks of size C by chopping up h on the time-axis with an overlap factor of 50%" → overlap 50% ⇒ stride of $C/2$

$$LN = Linear(h)_{F \times T'} \longrightarrow \boxed{Chunk} \longrightarrow h'' \in \mathbb{R}^{F \times C \times N_C}$$

$$Chunk(h;C): F \times T' \longmapsto F \times C \times N_C$$

Seq Frames

$h' \to$ [Intra Transformer] $\to$ [Repeat] $\to$ [Inter Transformer] $\to h''$

"Repeat N times"

$h'' = f_{inter}(P(f_{intra}(h')))$

Intra T ~ short term dependencies ~ applied second dim of $h'$
Inter T ~ long term dependencies ~ applied across chunks

P ~ permute last two dim.

Transformer Block with Residual Connections

$z \to \oplus \leftarrow c$, $\to z' \to$ [Layer Norm] $\to$ [MHA] $\to z'' \to \oplus \to$ [Layer Norm] $\to$ [FFW] $\to \oplus \to z''' \to \oplus \to f(z)$

"Repeat k times"

Intra and Inter Transformers

$z' = z + c$  // c ~ position encoding

$z'' = $ Multi Head Attention (Layer Norm $(z')$)

$z''' = $ Feed Forward (Layer Norm $(z'' + z')$) + $z'' + z'$

$f(z) = g^k(z + c) + z$

~ $g^k(\cdot)$ denotes k layers of transformers

" we add residual connections across the transformer ... across the transformers additive to improve gradient backprop "