

# Dual-Domain Joint Encoder 语音分离模型详解

## 1. 论文背景与模型目标

### 1.1 语音分离任务简介

语音分离（Speech Separation）是语音信号处理中的关键问题，旨在从混合语音信号中分离出各个独立的说话人。这在 **语音识别（ASR）、助听设备、会议转录** 等应用中至关重要。

传统的语音分离方法通常采用 **单一域信息**（仅时域或仅频域），但单独依赖某一域的信息可能存在局限性：

- 时域方法**（如 Conv-TasNet）能够直接处理波形信号，但难以捕捉语音的谐波和时频结构。
- 频域方法**（如 Deep Clustering、Spectral Masking）能够利用 STFT 变换后的频率特征，但会导致时序信息丢失。

### 1.2 论文核心思想：时域与频域联合编码

本论文提出 **Dual-Domain Joint Encoder** 语音分离模型，结合**时域特征**和**频域特征**，利用**联合编码和特征融合**的策略，充分利用语音的时频信息，提高语音分离质量。

核心贡献包括：

- 时域与频域联合编码（Dual-Domain Encoding）**：
  - 传统的语音分离方法仅利用时域或频域的信息，但单一域的信息通常无法充分表示复杂的语音混合信号。
  - 本模型同时在 **时域** 和 **频域** 进行特征提取，分别通过 **时域编码器（Time Domain Encoder）** 和 **频域编码器（Frequency Domain Feature Extractor）** 获取特征。
  - 结合 **卷积神经网络（CNN）** 进行特征提取，提高对语音信号的建模能力。
- 特征融合模块（Time-Frequency Fusion Module）**：
  - 通过 **领域融合（Domain Fusion）** 提取跨域的信息，使时域特征与频域特征能够互补。
  - 通过 **特征融合（Feature Fusion）** 进一步增强语音分离的判别能力，确保分离网络能够更准确地恢复目标语音信号。
- 自适应掩码预测（Mask Estimation）**：
  - 在频域上，通过学习 **掩码（Mask）**，提升信号的分离能力。
  - 掩码作用于 **STFT 变换的幅度谱**，然后通过 **逆 STFT（iSTFT）** 恢复时域信号。
  - 通过**端到端的训练**，模型能够直接优化时域的信号重建，而非仅优化频域的掩码预测。

### 1.3 论文方法的优势

- 跨域信息互补**：
  - 频域编码器可以从 STFT 变换的频谱中提取频率信息，而时域编码器则保留了语音的时序信息。
  - 通过 **特征融合模块（Fusion Module）**，两种特征能够有效结合，提高模型的表现力。
- 端到端训练**：
  - 传统的频域方法需要额外的 **后处理步骤（如 iSTFT）**，但本方法可以在训练过程中 **优化最终的波形质量**，从而减少伪影（artifacts）。
- 适用于多种语音环境**：
  - 该方法不仅可以用于 **无噪声环境**（clean speech separation），还可以扩展到 **带噪声的混合语音**，适用于更复杂的语音增强任务。

### 1.4 本文档目标

本文件详细解析 **Dual-Domain Joint Encoder 语音分离模型** 的代码实现，结合论文内容进行深入说明，涵盖以下内容：

- 模型结构**：时域编码、频域编码、特征融合、掩码估计
- 数据流动过程**：输入 → 编码 → 分离网络 → 掩码 → iSTFT 还原
- 关键超参数分析**：STFT 变换、卷积核大小、步长等
- 训练策略**：损失函数（SI-SNR）、优化器、批量大小等
- 实验结果及误差分析**

本模型的目标是 **复现论文结果**，并提供详细的实现说明，帮助研究人员和工程师理解 Dual-Domain Joint Encoder 在语音分离任务中的应用。

## 2. 模型结构

模型主要包含以下几个部分：

### 2.1 时域编码器（Time Domain Encoder）

#### 概述

在本模型中，时域编码器扮演着至关重要的角色，其主要任务是从原始音频波形中提取出能够描述信号内在结构的特征表示。时域编码器采用了 Conv-TasNet 的设计理念，通过一维卷积（1D Convolution）将原始音频信号转换到一个潜在空间中，从而使得后续的语音分离模块可以在这一高维、稀疏且具有判别能力的表示上进行操作。换句话说，每个卷积核在训练过程中都学习到了一种“基函数”，这些基函数构成了信号的字典，每个信号都可以看作是这些基函数的稀疏线性组合。

#### 数学模型与基本原理

一维卷积运算在信号处理和深度学习中被广泛应用，其核心思想是通过局部加权求和操作提取信号中的局部特征。设输入信号为  $\mathbf{x} \in \mathbb{R}^L$ （其中  $L$  为信号的长度），卷积核权重为  $\mathbf{w} \in \mathbb{R}^K$ （ $K$  为卷积核的大小），则卷积操作可以表示为：

$$y[t] = (\mathbf{x} * \mathbf{w})[t] = \sum_{k=0}^{K-1} w[k] \cdot x[t + k - P],$$

其中  $P$  表示填充（padding）参数， $t$  为输出信号的时间步。对于一个批量输入  $\mathbf{X} \in \mathbb{R}^{B \times 1 \times L}$ （ $B$  为批量大小，1 表示单通道），一维卷积将其映射到一个新的特征空间  $\mathbf{Y} \in \mathbb{R}^{B \times C \times L'}$ ，其中  $C$  为输出通道数，而  $L'$  则由卷积核大小、步长（stride）以及填充参数共同决定。

在 Conv-TasNet 的设计中，时域编码器通过一组无偏置的卷积核将输入信号投影到一个潜在表示空间中。数学上，如果我们设定卷积核的数量为  $N$ （即输出通道数 encoder\_dim），那么对于每个滤波器  $j = 1, 2, \dots, N$ ，对应的卷积运算为：

$$y_j[t] = (\mathbf{x} * \mathbf{w}_j)[t] = \sum_{i=0}^{K-1} w_j[i] \cdot x[t + i - P],$$

将所有  $N$  个特征图合并，我们获得的特征表示  $\mathbf{Y} \in \mathbb{R}^{N \times L'}$  就是该时域编码器的输出。这种表示能够捕捉到信号在不同时间尺度上的局部结构，进而为后续的语音分离提供有力的特征支持。

#### 参数设置的详细分析

在代码中，时域编码器的定义如下：

```
self.time_encoder = nn.Conv1d(1, encoder_dim, kernel_size=16, stride=8, bias=False)
```

下面逐一对各个参数进行详细说明：

- 1. 输入通道数（in\_channels = 1）**  
这里设置为1，意味着输入的音频信号是单通道信号。在大多数语音处理任务中，我们通常处理单声道音频，这样可以直接用一个通道来描述信号。如果处理多通道音频（例如立体声），则需要将输入通道数设为对应的通道数。
- 2. 输出通道数（out\_channels = encoder\_dim）**  
参数 encoder\_dim 通常设定为 256，表示时域编码器将输出 256 个不同的特征图。这些特征图可以视为 256 个基函数，每个基函数在训练过程中会自动学习到能够捕捉输入信号特定模式的形状。增加输出通道数通常会提升模型对复杂信号的捕捉能力，但同时也会增加计算复杂度和模型参数量。
- 3. 卷积核大小（kernel\_size = 16）**  
卷积核大小是指每次卷积操作考虑的信号采样点数。在本例中，卷积核大小设为 16，意味着每个滤波器在做局部加权求和时，会同时考虑 16 个连续的采样点。数学上，假设滤波器  $\mathbf{w}$  的长度为 16，则在位置  $t$  上的卷积操作为：

$$y[t] = \sum_{i=0}^{15} w[i] \cdot x[t + i - P].$$

这种设置能够捕捉到信号中相对较短时长内的变化，比如瞬时频率变化或微小的波形细节，且有助于后续稀疏表示的构建。

- 4. 步长（stride = 8）**  
步长决定了卷积核在输入信号上滑动的间隔。这里设置步长为 8 表示卷积核每次向前移动 8 个采样点，进而对信号进行下采样。数学上，若原始信号长

度为  $L$ ，那么经过卷积操作后的输出长度  $L'$  可以通过下式计算：

ParseError: KaTeX parse error: Expected 'EOF', got '\_' at position 41: ... - \text{kernel\\_size} + 2P\}\te...

当  $P$  为 0 或经过合理的填充后，步长的设定会显著降低时间维度的长度，同时保留主要特征。下采样的好处在于降低后续处理的计算量，同时在一定程度上使得信号表示更具紧凑性。

5. 偏置 (bias = False)

在本设计中，时域编码器的卷积层不使用偏置项。这是因为在 Conv-TasNet 的设计中，我们希望卷积核仅代表输入信号中的基函数，而不受额外常数偏置的影响。数学上去除偏置项可以减少参数数量，并使得卷积运算仅依赖于输入信号与滤波器的内积，从而更直接地反映信号的局部结构。

设计理念与 Conv-TasNet 的关联

Conv-TasNet 模型提出了一种基于时域的语音分离方法，其核心思想在于直接在原始时域上进行操作，而非依赖传统的短时傅里叶变换（STFT）。这种方法的优势在于避免了传统频域方法中常见的相位重构问题，同时通过学习到的基函数实现了信号的稀疏表示。在 Conv-TasNet 中，时域编码器与其对应的解码器构成了一个编码-解码框架，编码器将原始信号投影到一个潜在表示空间，而解码器则利用转置卷积将这一表示还原为时域信号。

数学上，我们可以将时域编码器视为一个线性变换，其作用为：

$$\mathbf{Y} = \mathcal{E}(\mathbf{x}) = \mathbf{W} * \mathbf{x},$$

其中  $\mathbf{W} \in \mathbb{R}^{encoder\_dim \times 1 \times kernel\_size}$  表示所有卷积核构成的权重矩阵，符号  $*$  表示一维卷积运算。这种变换将原始信号  $\mathbf{x}$  映射到一个高维空间  $\mathbf{Y} \in \mathbb{R}^{encoder\_dim \times L'}$  中，每个通道均对应一个可学习的基函数。通过这种方式，模型能够在训练过程中自动优化这些基函数，使其具有区分不同语音信号的能力，并在语音分离任务中提供充分的特征支持。

卷积运算中的细节与实现考虑

在实际实现中，卷积运算的效率和稳定性都至关重要。首先，在 PyTorch 框架下，nn.Conv1d 提供了高效的卷积实现，其内部会自动利用底层的高性能库（如 cuDNN）加速运算。其次，在设计时域编码器时，需特别注意输出时间步数的计算，因为这会影响到后续模块（例如频域分支、融合模块和解码器）的设计。设输入信号长度为  $L$ ，经过无偏置卷积后的输出长度为：

$$L' = \left\lfloor \frac{L - kernel\_size + 2P}{stride} \right\rfloor + 1.$$

合理选择 kernel\_size、stride 和 padding 参数，可以确保不同分支的特征图在时间维度上对齐，从而便于信息融合。

此外，由于时域编码器的输出将直接影响到后续的掩码预测与语音重构，如何确保该模块提取到的信息具有足够的判别性和稳定性，成为设计的重要考量。为此，在 Conv-TasNet 中通常会在训练过程中采用一定的正则化手段（如 dropout 或权重衰减）来防止过拟合，同时采用合适的优化算法（如 Adam）加速收敛。所有这些设计均围绕一个核心目标：使得时域编码器能以最优方式将原始信号映射到潜在空间，捕捉到足够丰富的局部特征信息。

与传统方法的对比

传统的语音分离方法往往依赖于 STFT 等频域方法，虽然能有效分离信号，但在重构阶段常常面临相位恢复困难的问题。而 Conv-TasNet 的时域编码器直接对原始信号进行处理，避免了频域方法中由于相位信息缺失而引起的重构误差。数学上，这种方法实现了从时域到潜在空间的直接线性映射，使得模型可以在高维空间中进行更加灵活和精细的信号分离。相比之下，传统方法需要先将信号转换为幅度谱，再利用掩码进行分离，最后再进行逆变换，这个过程不仅复杂，而且由于相位信息的处理不当，容易导致分离后的信号失真。

训练中的数值优化

在训练阶段，时域编码器的权重  $\mathbf{W}$  会随着损失函数的反向传播不断更新。设定损失函数为  $\mathcal{L}$ ，对权重  $\mathbf{W}$  的梯度为：

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \sum_t \frac{\partial \mathcal{L}}{\partial y[t]} \cdot \frac{\partial y[t]}{\partial \mathbf{W}},$$

这一公式表明，每个卷积核在整个训练数据上的误差信息都会累积，进而引导模型不断调整滤波器的形状。随着训练的进行，这些滤波器逐渐学习到能够准确捕捉语音信号特定结构的基函数，这种学习过程正是 Conv-TasNet 成功的关键所在。

总结

时域编码器是 Conv-TasNet 设计中的核心模块，其主要任务是通过一维卷积操作将原始音频信号映射到一个潜在的高维表示空间。在这一过程中，每个卷积核都扮演着学习到的基函数的角色，这些基函数共同构成了信号的字典。通过合理的参数设置（如 kernel\_size=16、stride=8 以及不使用偏置），时域编码器不仅能够捕捉信号中丰富的局部特征，而且还能对信号进行有效的下采样，为后续的语音分离和重构提供紧凑且判别性强的特征表示。与传统依赖频域转换的方法相比，时域编码器直接处理原始信号，避免了相位恢复等问题，从而在实际应用中展示了更高的鲁棒性和效果。

2.2 频域特征提取（Frequency Domain Feature Extractor）

概述

在语音处理任务中，频域分析一直被认为是揭示信号内在结构的重要工具。传统的语音分离方法常借助短时傅里叶变换（STFT）将一维时域信号转换为二维的时频表示，这一表示不仅能够捕捉到信号在时间上的变化，还能反映出频谱的分布情况。频域表示提供了丰富的语音特征信息，如共振峰、噪声分布和语音的谐波结构等。频域特征提取模块的核心思想是，首先利用 STFT 将原始波形  $\mathbf{x}(t)$  转换为频谱表示  $\mathbf{X}(f, \tau)$ ，其中  $f$  表示频率 bin， $\tau$  表示时间帧，然后利用一维卷积层对频谱进行进一步处理，以提取更具判别性的频域特征。这一过程不仅强化了模型对局部频谱模式的捕捉能力，同时也为后续的时频融合提供了稳健的特征输入。

短时傅里叶变换（STFT）数学描述

短时傅里叶变换是一种将时域信号分解为时间和频率双重域的方法。设原始时域信号为  $\mathbf{x}(t)$ ，STFT 定义为：

$$X(\tau, \omega) = \sum_{t=-\infty}^{\infty} x(t)w(t - \tau)e^{-j\omega t},$$

其中：

- $\tau$  表示窗函数  $w(t - \tau)$  的中心位置，对应时间帧；
- $\omega$  是角频率；
- $w(t)$  为窗函数，常见的窗函数包括汉明窗、汉宁窗或矩形窗等，用于限制信号在局部区域内的计算，从而获得局部频谱特征。

在离散时间的实现中，我们将信号  $\mathbf{x}[n]$  分成重叠的帧，每一帧通过窗函数后进行傅里叶变换。设采样率为  $f_s$ ，每一帧长度为  $L$ （通常等于 FFT 长度  $n_{\text{fft}}$ ），则对于每个时间帧  $\tau$  对应的频谱可以写成：

$$X[\tau, k] = \sum_{n=0}^{L-1} x[n + \tau H] w[n] e^{-j \frac{2\pi}{L} kn},$$

其中  $H$  为帧移（hop length）， $k = 0, 1, \dots, L - 1$  表示频率 bin。在实际计算中，为了减少冗余信息，通常只保留非负频率部分，因此频谱的维度为  $\frac{n_m}{2} + 1$ 。

对于本模型来说，STFT 提供了一个二维表示，其形状通常为  $(B, \text{freq\_bins}, T)$ ，其中  $B$  表示批量大小， $\text{freq\_bins} = \frac{n_m}{2} + 1$  表示频率 bin 数，而  $T$  为时间帧数。STFT 的幅度谱  $|X(\tau, k)|$  具有很强的鲁棒性，并能反映出语音信号中的共振峰和能量分布，因而在语音分离任务中起到关键作用。

1D 卷积在频域特征提取中的作用

在将原始时域信号经过 STFT 转换到频域后，我们得到一个复数矩阵  $\mathbf{X}$ ，但实际应用中常对其取幅度谱  $|\mathbf{X}|$  作为输入特征。此时，特征矩阵的形状为  $(B, \text{freq\_bins}, T)$ 。为了进一步提取局部频域模式，本模型采用了一维卷积层来对幅度谱进行特征提取。具体而言，模型使用了如下定义的卷积层：

```
self.freq_encoder = nn.Conv1d(n_fft // 2 + 1, encoder_dim, kernel_size=3, padding=1)
```

这里的 `nn.Conv1d` 接受的输入通道数为  $\frac{n_m}{2} + 1$ ，即每个频率 bin 作为一个独立的输入通道。这种做法在数学上等价于将频谱矩阵视为具有多个通道的一维信号，其中通道数对应于频率分量。

在卷积操作中，设输入特征  $\mathbf{F} \in \mathbb{R}^{C_{\text{in}} \times T}$ （这里  $C_{\text{in}} = \frac{n_m}{2} + 1$ ），卷积核  $\mathbf{K} \in \mathbb{R}^{C_{\text{in}} \times K}$ （ $K = 3$  为卷积核大小），则卷积操作输出的特征在第  $j$  个通道上可以表示为：

$$y_j[t] = \sum_{i=0}^{C_{\text{in}}-1} \sum_{n=0}^{K-1} K_{j,i,n} \cdot F[i, t + n - P],$$

其中  $P = 1$  是填充长度，保证输出特征的时间步数与输入保持一致。这一操作能够捕捉到局部时间窗口内不同频率通道之间的相关性，同时对每个时间步的频谱信息进行组合，从而提取出更高级别的频域特征表示。输出的特征维度为  $(B, \text{encoder\_dim}, T)$ ，其中 `encoder_dim` 为设定的输出通道数（通常与时域编码器保持一致，便于后续的特征融合）。

参数选择与设计细节

1. **输入通道数**  $n_{\text{fft}} // 2 + 1$   
STFT 计算时，仅保留非负频率部分，使得输入通道数为  $\frac{n_{\text{fft}}}{2} + 1$ 。这在数学上可以解释为利用对称性原理：对于实值信号，其傅里叶变换满足共轭对称性，因此只保留前半部分即可。这不仅降低了计算复杂度，同时保留了所有有用的频域信息。
2. **卷积核大小**  $kernel\_size = 3$   
选择较小的卷积核大小（如 3）是为了在局部范围内捕捉细粒度的频域变化。数学上，这相当于在每个时间步上对相邻 3 个频率 bin 进行局部加权求和，这种方式能够有效捕捉到频谱中细微但具有辨识度的模式，如谐波结构或噪声分布的局部变化。较小的卷积核能够降低参数量，防止过拟合，同时使得网络具有更好的泛化能力。
3. **填充**  $padding = 1$   
设置填充大小为 1 是为了保持卷积操作前后特征图的时间步数一致。数学上，通过在输入信号两侧补 1 个零，可以保证输出长度为：

$$T_{\text{out}} = \left\lfloor \frac{T + 2P - K}{\text{stride}} \right\rfloor + 1 = T,$$

其中  $K = 3$  且  $P = 1$  且 stride 默认为 1。保持输入和输出维度一致有助于后续的特征融合模块直接将频域特征与时域特征拼接，而无需额外调整时间步长。

频域特征提取的优势与作用

在语音分离任务中，频域特征提取具有以下几个显著优势：

1. **捕捉局部频谱结构**  
语音信号在频域中常展现出明显的谐波结构和共振峰信息。通过 STFT 得到的幅度谱提供了关于能量分布和频率分量的详细描述，而一维卷积则能够在每个时间步内捕捉这些局部模式。这种局部特征不仅能区分不同的语音源，还能在噪声干扰较强的情况下提取出有用信息。
2. **降维与特征压缩**  
直接处理原始的 STFT 幅度谱可能会面临维度较高和冗余信息较多的问题。通过卷积层进行特征提取，可以将高维频域信息映射到一个更紧凑的特征空间中，这样不仅降低了计算复杂度，还能过滤掉噪声和不必要的冗余信息，增强模型的鲁棒性。
3. **便于与时域特征融合**  
在多域联合建模中，将时域和频域特征进行融合是提升分离效果的重要策略。通过卷积提取后的频域特征，其形状与时域编码器输出的特征具有一致的时间步长和相近的维度，这使得后续通过拼接和融合操作能够实现信息互补，从而为语音分离提供更全面的特征描述。
4. **局部平移不变性**  
卷积操作本身具有局部平移不变性，这意味着即使语音信号在时间轴上存在轻微的偏移，其局部频谱结构仍能被正确捕捉。这一特性在实际应用中尤为重要，因为语音信号通常会受到说话人速度变化、环境噪声等因素的干扰，而局部平移不变性则能增强模型对这些变化的鲁棒性。

实现细节与训练注意事项

在实际实现过程中，频域特征提取模块通常与 STFT 操作紧密结合。具体流程为：

1. 对输入信号  $\mathbf{x}(t)$  应用 STFT，得到复数频谱  $\mathbf{X}(\tau, k)$ 。
2. 对  $\mathbf{X}$  取幅度  $|\mathbf{X}|$  并扩展维度，使其形状为  $(B, 1, \text{freq\_bins}, T)$ 。
3. 使用 `squeeze(1)` 操作后，将数据形状转换为  $(B, \text{freq\_bins}, T)$ ，此时每个频率 bin 作为一个通道输入到 Conv1d 层中。
4. 卷积层对输入数据进行滑动窗口操作，利用局部加权求和提取出频域特征，输出的特征维度为  $(B, \text{encoder\_dim}, T)$ 。

在训练过程中，频域特征提取模块的卷积核权重会随着反向传播不断更新。损失函数会根据分离结果与目标语音之间的差异计算梯度，从而指导卷积核参数的调整。通过这种方式，卷积核能够逐渐学习到如何捕捉频谱中的关键特征，如语音的谐波结构和噪声模式，从而提升整体模型的分离性能。

此外，合理的初始化方法和正则化策略也对频域特征提取模块的训练效果至关重要。通常会采用 Xavier 或 He 初始化方法为卷积核赋初值，以确保梯度在网络中传播时不会消失或爆炸。同时，可能结合 Dropout 或权重衰减等手段防止模型过拟合，提高模型在实际应用中的泛化能力。

总结

频域特征提取模块利用 STFT 将时域信号转换为时频表示，再通过一维卷积对幅度谱进行局部特征提取。这一过程从数学上将语音信号映射到一个低维且具有丰富判别信息的特征空间中，为后续时频融合提供了坚实的基础。具体来说，通过设定输入通道数为  $\frac{n_{\text{fft}}}{2} + 1$ （对应 STFT 计算的频率 bin 数）、采用  $\text{kernel\_size} = 3$  的小卷积核以及  $\text{padding} = 1$  保证输出尺寸一致，频域特征提取模块在捕捉局部频谱模式、压缩冗余信息以及增强模型鲁棒性方面均发挥了重要作用。

2.3 时频融合模块（Time-Frequency Fusion）

概述

时频融合模块在多域语音分离系统中起到至关重要的作用。由于语音信号在时域和频域中都包含有重要的互补信息，如何有效地融合这两种特征以获得更具判别性的表示，是提高分离性能的关键。时域特征往往能够捕捉到信号的瞬时变化和局部波形结构，而频域特征则反映了信号的谐波成分、共振峰以及频谱分布。二者各有所长，相互融合后能够为后续的分离网络提供更丰富的上下文信息。时频融合模块主要依靠点卷积（即卷积核大小为 1 的 Conv1d）来实现这一信息整合，同时引入非线性激活函数 ReLU 以增强模型表达能力。

数学描述与基本原理

假设我们已经分别通过时域编码器和频域特征提取模块获得了两个特征表示：

- 时域特征表示： $\mathbf{F}_{\text{time}} \in \mathbb{R}^{B \times C \times T}$
- 频域特征表示： $\mathbf{F}_{\text{freq}} \in \mathbb{R}^{B \times C \times T}$

其中， $B$  表示批量大小， $C$  表示编码器输出的通道数（即 `encoder_dim`）， $T$  表示时间步数。通常，这两个特征在时间步数上已经通过插值等操作进行了对齐，保证它们在时间维度上可以直接拼接。

时频融合模块的第一步是将这两个特征在通道维度上进行拼接，得到新的特征表示：

$$\mathbf{F}_{\text{concat}} = \text{Concat}(\mathbf{F}_{\text{time}}, \mathbf{F}_{\text{freq}}) \in \mathbb{R}^{B \times 2C \times T}$$

拼接操作在数学上可以看作是将两个矩阵沿着通道维度进行堆叠，即对于每个时间步  $t$  有：

$$\mathbf{f}_{\text{concat}}(t) = \begin{bmatrix} \mathbf{f}_{\text{time}}(t) \\ \mathbf{f}_{\text{freq}}(t) \end{bmatrix} \in \mathbb{R}^{2C}$$

这一步骤实现了将时域和频域信息进行融合，为后续的非线性变换提供了丰富的信息基础。

随后，融合后的特征会经过一个由两个点卷积层（`kernel_size=1`）和中间的 ReLU 激活函数组成的顺序模块。该模块的结构如下：

```
self.domain_fusion = nn.Sequential(  
    nn.Conv1d(encoder_dim * 2, encoder_dim, kernel_size=1),  
    nn.ReLU(),  
    nn.Conv1d(encoder_dim, encoder_dim, kernel_size=1)  
)
```

我们可以对每一层的数学运算进行详细分析。

第一层点卷积

第一层卷积操作的输入是  $\mathbf{F}_{\text{concat}} \in \mathbb{R}^{B \times 2C \times T}$ ，使用卷积核大小为 1、输入通道数  $2C$  和输出通道数  $C$  进行变换。记该卷积核为  $\mathbf{W}_1 \in \mathbb{R}^{C \times 2C \times 1}$ ，则对于每个时间步  $t$ ，输出特征  $\mathbf{z}(t) \in \mathbb{R}^C$  的计算公式为：

$$z_i(t) = \sum_{j=1}^{2C} W_{1,i,j} \cdot f_{\text{concat},j}(t) + b_{1,i}, \quad i = 1, 2, \dots, C,$$

其中  $b_{1,i}$  是第  $i$  个输出通道的偏置项（在许多实现中也可以选择不使用偏置项）。这里的卷积核大小为 1 表明，卷积操作只在通道维度上做线性组合，而不涉及时间步之间的信息混合。该操作的核心思想是对拼接后的高维特征进行线性变换，降维到原始通道数  $C$ ，同时学习如何从时域和频域信息中提取出更有判别力的特征组合。

激活函数 ReLU

经过第一层点卷积后，特征  $\mathbf{z}(t)$  会传递给 ReLU 激活函数：

$$\mathbf{a}(t) = \text{ReLU}(\mathbf{z}(t)) = \max(0, \mathbf{z}(t)).$$

ReLU 的作用在于引入非线性变换，能够有效破除线性变换的局限，使得模型可以表达更复杂的函数映射。数学上，ReLU 将负数截断为零，从而稀疏化激活值，并在一定程度上缓解梯度消失问题。这一步骤确保了时频融合模块不仅仅停留在线性变换，而是能够捕捉到时频信息之间的非线性关系。

第二层点卷积

在经过激活函数后，融合模块接下来应用第二个点卷积层。设该卷积核为  $\mathbf{W}_2 \in \mathbb{R}^{C \times C \times 1}$ ，那么对于每个时间步  $t$ ，第二层卷积的输出  $\mathbf{y}(t) \in \mathbb{R}^C$  可以表示为：

$$y_i(t) = \sum_{j=1}^C W_{2,i,j} \cdot a_j(t) + b_{2,i}, \quad i = 1, 2, \dots, C.$$

该操作进一步对已经融合且经过非线性激活的特征进行线性变换，旨在提炼出最终用于后续分离网络处理的高质量特征。点卷积的特点在于其参数量较少，计算效率高，并且能够实现逐通道的信息重新加权 and 组合。

融合模块整体效果与优势

综合来看，时频融合模块利用点卷积将拼接后的时域和频域特征进行降维和信息重组，其核心思想可以归纳为：

- 1. **信息拼接**：通过对时域和频域特征进行拼接，将两个域的特征堆叠到一起，形成一个更为丰富的特征表示。这种方法在数学上等价于直接将两个向量连接起来，保持各自的信息完整性。
- 2. **降维与重组**：第一层点卷积将拼接后的高维特征降到原始的  $C$  维空间，同时学习到如何在两种不同域的信息中提取有用的特征。随后经过 ReLU 的非线性映射，再经过第二层点卷积进一步整合信息，输出的特征既包含时域信息的局部动态，又保留了频域信息的谐波和共振模式。
- 3. **非线性映射**：ReLU 激活函数的引入使得融合过程具备非线性表达能力，能够捕捉到复杂的特征关系，而不仅仅局限于线性组合。非线性映射能够增强模型的表达能力和判别能力，从而更好地适应语音分离等复杂任务。

这种设计的优势在于，虽然点卷积本身只是对通道进行线性组合，但结合非线性激活后，整个模块能够实现类似于全连接层的功能，同时保留时序数据的局部结构不变性。由于卷积核大小为 1，因此计算复杂度相对较低，适合大规模数据处理，并且能够与其他模块（例如时域编码器和频域提取模块）无缝对接，保证整个网络在前向传播中的维度一致性。

参数设置与实现考虑

- **输入通道数**  $encoder\_dim \times 2$   
拼接时域和频域特征后，通道数由原始的  $encoder\_dim$  翻倍至  $2 \times encoder\_dim$ 。这种设计确保了两种域的信息能够在融合前都得到完整保留，为后续降维和特征重组提供充足的信息冗余。
- **卷积核大小**  $kernel\_size = 1$   
使用点卷积的核心原因在于其仅在通道维度上进行运算，而不改变时间步数。数学上，点卷积对每个时间步  $t$  都独立地应用相同的线性变换，这使得时间信息保持不变，同时允许不同通道间进行信息融合和交互。
- **ReLU 激活函数**  
激活函数在深度神经网络中扮演着非线性映射的重要角色，ReLU 简单且高效，其非线性特性能够使得融合后的特征具备更强的表达能力和鲁棒性。ReLU 的导数在正区间恒为 1，有助于缓解梯度消失问题，加速网络训练。



• 连续两层卷积结构

两层连续的点卷积（中间夹有 ReLU）形成了一个轻量级的全连接网络结构，其参数量远低于传统全连接层，但却能有效实现特征维度的转换和重组。这种设计既保持了计算效率，又能充分发挥融合操作的潜力。

在整体网络中的作用

时频融合模块位于时域编码器和频域提取模块之后，以及分离网络之前，其输出将直接影响到后续语音分离的效果。通过融合操作，模型能够在一个统一的特征空间中同时利用时域和频域信息，从而提高对不同语音源特征的分辨能力。具体来说，融合模块输出的特征不仅包含了时域上捕捉到的瞬时能量和波形变化信息，还蕴含了频域上展现的频谱分布、共振峰和谐波结构。这样的信息整合使得后续的分离网络能够更准确地计算每个语音源的掩码，进而实现更高质量的语音重构。

从数学上讲，时频融合模块实现了一个从高维联合特征空间到低维判别特征空间的映射，即：

$$\mathbf{F}_{\text{fused}} = \phi(\text{Concat}(\mathbf{F}_{\text{time}}, \mathbf{F}_{\text{freq}})),$$

其中  $\phi$  代表由点卷积和非线性激活构成的映射函数。这种映射保证了信息在维度降低的同时，其判别性和表达能力得到强化，能够为语音分离任务提供充分的支持。

2.4 分离网络（Separation Network）

1. 模块概述

分离网络主要由两层 1D 卷积层和对应的 ReLU 激活函数构成。其核心目标在于对融合后的时频特征进行进一步提取和非线性映射，从而学习到更具判别性的特征表示，为后续的语音掩码预测提供高质量的输入。该模块采用的结构为：

```
self.separation_net = nn.Sequential(  
    nn.Conv1d(encoder_dim, encoder_dim, kernel_size=3, padding=1),  
    nn.ReLU(),  
    nn.Conv1d(encoder_dim, encoder_dim, kernel_size=3, padding=1),  
    nn.ReLU()  
)
```

在这一模块中，输入和输出的通道数均为 `encoder_dim`，并且两个卷积层均采用了 `kernel_size=3` 和 `padding=1`，确保特征图的时间步数在卷积操作前后保持一致。

2. 数学原理与卷积运算细节

2.1 卷积操作的基本定义

对于一维卷积，设输入特征图为  $\mathbf{X} \in \mathbb{R}^{C \times T}$ （这里  $C = \text{encoder\_dim}$  为通道数， $T$  为时间步数），卷积核为  $\mathbf{K} \in \mathbb{R}^{C \times k}$ ，其中  $k = 3$  是卷积核大小。对于每个输出通道上的时间步  $t$ ，卷积运算可以表示为：

$$y(t) = \sum_{c=1}^C \sum_{i=0}^{k-1} K(c, i) \cdot X(c, t + i - p) + b,$$

其中  $p$  为填充（padding）大小， $b$  是偏置项（本例中没有额外说明是否使用偏置，但一般卷积层默认包含偏置项，除非显式设置为 `bias=False`）。在本例中，使用了 `padding=1`，意味着在时间轴两侧各补 1 个零，使得输出的时间步数  $T_{\text{out}}$  与输入  $T$  保持不变，即：

$$T_{\text{out}} = \left\lfloor \frac{T + 2p - k}{\text{stride}} \right\rfloor + 1 = \left\lfloor \frac{T + 2 \times 1 - 3}{1} \right\rfloor + 1 = T.$$

2.2 第一层卷积与 ReLU 激活

第一层卷积

该层采用了输入和输出均为 `encoder_dim` 的卷积层，其数学计算公式为：



$$\mathbf{Z}^{(1)}(t) = \mathbf{W}^{(1)} * \mathbf{F}(t) + \mathbf{b}^{(1)},$$

其中：

- $\mathbf{F}(t) \in \mathbb{R}^{encoder\_dim}$  表示融合后的特征在时间步  $t$  的向量表示，
- $\mathbf{W}^{(1)} \in \mathbb{R}^{encoder\_dim \times encoder\_dim \times 3}$  是第一层卷积核，
- $\mathbf{b}^{(1)} \in \mathbb{R}^{encoder\_dim}$  为偏置项，
- “\*” 表示一维卷积运算。

由于采用了 `padding=1`，对于每个时间步  $t$ ，卷积核会覆盖  $t - 1$ 、 $t$  和  $t + 1$  三个相邻时刻的特征，从而捕捉局部的时间依赖关系。用数学语言描述，这意味着：

$$z_i^{(1)}(t) = \sum_{j=1}^{encoder\_dim} \left[ w_{i,j}^{(1)}(-1) \cdot f_j(t-1) + w_{i,j}^{(1)}(0) \cdot f_j(t) + w_{i,j}^{(1)}(1) \cdot f_j(t+1) \right] + b_i^{(1)},$$

其中  $w_{i,j}^{(1)}(k)$  表示第  $i$  个输出通道对应第  $j$  个输入通道在卷积核中第  $k$  个权重。

**ReLU 激活**

接着，卷积后的输出通过 ReLU 非线性激活函数：

$$a_i^{(1)}(t) = \text{ReLU} \left( z_i^{(1)}(t) \right) = \max \left( 0, z_i^{(1)}(t) \right).$$

ReLU 的作用在于引入非线性变换，使得模型能够拟合更加复杂的函数，同时抑制负值的传播，起到稀疏化激活的效果。

**2.3 第二层卷积与 ReLU 激活**

**第二层卷积**

将第一层经过 ReLU 激活后的输出记为  $\mathbf{A}^{(1)}(t)$ ，第二层卷积层接收该输出并进行类似的运算：

$$\mathbf{Z}^{(2)}(t) = \mathbf{W}^{(2)} * \mathbf{A}^{(1)}(t) + \mathbf{b}^{(2)},$$

其中  $\mathbf{W}^{(2)} \in \mathbb{R}^{encoder\_dim \times encoder\_dim \times 3}$  是第二层卷积核， $\mathbf{b}^{(2)} \in \mathbb{R}^{encoder\_dim}$  为偏置项。同样地，填充  $p = 1$  保证了第二层卷积的输出在时间维度上与输入保持一致。

数学上，对第二层卷积的每个输出特征  $z_i^{(2)}(t)$  可以表达为：

$$z_i^{(2)}(t) = \sum_{j=1}^{encoder\_dim} \left[ w_{i,j}^{(2)}(-1) \cdot a_j^{(1)}(t-1) + w_{i,j}^{(2)}(0) \cdot a_j^{(1)}(t) + w_{i,j}^{(2)}(1) \cdot a_j^{(1)}(t+1) \right] + b_i^{(2)}.$$

**第二个 ReLU 激活**

随后，再次应用 ReLU 激活函数：

$$a_i^{(2)}(t) = \text{ReLU} \left( z_i^{(2)}(t) \right) = \max \left( 0, z_i^{(2)}(t) \right).$$

这一激活过程进一步引入非线性特征，并抑制了那些不具备判别信息的负激活值，使得最终的输出特征  $\mathbf{A}^{(2)}(t)$  能够更准确地表达各个时间步内局部特征之间的复杂关系。

**3. 代码解析与实现细节**

在代码实现上，分离网络模块被封装在 `nn.Sequential` 中，从而形成一个顺序执行的子网络。具体代码如下：

```
self.separation_net = nn.Sequential(  
    nn.Conv1d(encoder_dim, encoder_dim, kernel_size=3, padding=1),  
    nn.ReLU(),  
    nn.Conv1d(encoder_dim, encoder_dim, kernel_size=3, padding=1),  
    nn.ReLU()  
)
```

- **nn.Conv1d(encoder\_dim, encoder\_dim, kernel\_size=3, padding=1)**  
这行代码实现了第一层一维卷积，输入和输出通道数均为 `encoder_dim`。设置 `kernel_size=3` 表示在每个时间步上，卷积核覆盖当前及相邻左右各一个时刻的特征。`padding=1` 则通过在输入的前后各补 1 个零来确保输出的时间步数不因卷积而减少。
- **nn.ReLU()**  
紧接着的 `nn.ReLU()` 模块应用了 ReLU 激活函数，将第一层卷积后的输出中所有负值置零，并保持正值不变。
- **第二个卷积和 ReLU**  
紧随其后，第二个 `nn.Conv1d` 与 `nn.ReLU()` 组合重复了上述过程，对经过第一层处理后的特征进一步进行局部信息融合和非线性变换。这两层组合形成了一个简洁且高效的特征提取器，其主要目的是在保证时间步数不变的前提下，对输入特征进行更深层次的抽象和重构。

4. 数学优势与局部时间依赖捕捉

采用 `kernel_size=3` 的卷积核能够非常有效地捕捉局部时间依赖关系。原因如下：

1. **局部依赖建模**  
对于语音信号，其局部结构往往包含重要的瞬时信息。例如，音素边界、瞬时共振峰变化等，这些信息通常在相邻三个时间步内存在较强的相关性。通过卷积操作，每个输出单元  $a_i^{(1)}(t)$  或  $a_i^{(2)}(t)$  都能够综合前一时刻、当前时刻和后一时刻的特征，从而形成更稳定和具有判别性的特征描述。
2. **参数共享与计算效率**  
卷积核在时间轴上进行滑动操作，参数共享的机制大大减少了参数数量，同时保持了模型对局部特征的灵活捕捉能力。数学上，这种共享机制使得卷积层可以看作是对整个序列应用相同的局部滤波器，确保在不同时间步上具有一致的特征提取效果。
3. **保持特征图尺寸**  
由于填充  $padding = 1$  的使用，卷积操作不会改变时间步数，这对于后续特征拼接或融合模块尤为重要。保持尺寸不变的设计使得后续网络在处理多尺度特征时更加方便，减少了因尺寸不匹配而引起的调整操作，保证了数据流在各模块之间的顺畅传递。
4. **非线性激活的作用**  
通过 ReLU 激活函数的引入，卷积层不仅实现了线性特征组合，还能够引入非线性因素，使得模型能够拟合更复杂的特征关系。非线性激活可以看作是在原始线性组合基础上添加了一个门控机制，只有当局部特征满足一定条件时，信息才会被进一步传递，这种机制在捕捉语音信号中的非线性模式方面起到了关键作用。

5. 模块在整体模型中的作用

在整个语音分离系统中，分离网络位于时频融合模块之后，其输入为已经融合了时域和频域信息的特征表示  $\mathbf{F}_{\text{fused}} \in \mathbb{R}^{B \times \text{encoder\_dim} \times T}$ 。分离网络通过两层卷积进一步提炼和加工这一表示，目标是提取出对语音分离任务最有利的特征，最终为语音掩码的预测提供基础。数学上，这可以表示为：

$$\mathbf{F}_{\text{sep}} = \psi(\mathbf{F}_{\text{fused}}),$$

其中  $\psi(\cdot)$  代表由分离网络实现的非线性映射函数，该函数由两层卷积和 ReLU 激活构成。通过这种映射，模型能够捕捉到输入特征中更高阶的局部相关性和非线性特征，使得后续的掩码预测模块在对各语音源进行分离时具有更高的准确性和鲁棒性。

2.5 掩码预测（Mask Estimation）

1. 模块概述

在语音分离任务中，掩码预测模块的主要目标是为每个时间步生成一个“掩码”（mask），该掩码用于对输入信号的频谱或时域特征进行加权，从而分离出不同说话人的语音。这里的“掩码”可以理解为一个在  $[0, 1]$  范围内的软门控信号，它决定了在每个时间步、每个频率点上各个声源的贡献。对于本模型而言，由于目标是分离两个说话人，故最终需要预测出 2 个独立的掩码。

该模块通过一个一维卷积层实现，其核心代码如下：

```
self.mask_predictor = nn.Conv1d(encoder_dim, 2, kernel_size=1) # 预测两个分离掩码
```

这行代码表示，将输入的特征（维度为 `encoder_dim`）在每个时间步上经过一个核大小为 1 的卷积操作，直接输出 2 个通道，对应于两个分离语音的掩码。接下来将详细说明这一过程的数学原理和实现细节。

2. 数学描述与卷积操作

2.1 核心数学公式

假设经过分离网络后的特征为

$$\mathbf{F} \in \mathbb{R}^{B \times C \times T},$$

其中：

- $B$  为批量大小，
- $C = \text{encoder\_dim}$  表示特征通道数，
- $T$  表示时间步数。

在掩码预测模块中，我们希望将每个时间步的  $C$  维特征向量映射到一个 2 维向量，分别对应两个说话人的掩码值。对于每个时间步  $t$  ( $t = 1, 2, \dots, T$ ) 以及批量中的每个样本，其特征向量记为：

$$\mathbf{f}(t) = [f_1(t), f_2(t), \dots, f_C(t)]^\top \in \mathbb{R}^C.$$

掩码预测操作可表示为一个线性变换：

$$\mathbf{m}(t) = \mathbf{W} \mathbf{f}(t) + \mathbf{b} \in \mathbb{R}^2,$$

其中：

- $\mathbf{W} \in \mathbb{R}^{2 \times C}$  为权重矩阵，
- $\mathbf{b} \in \mathbb{R}^2$  为偏置向量。

由于卷积核大小设置为 1，所以这一步操作在每个时间步  $t$  上独立进行，并且不会混合不同时间步之间的信息。也就是说，1D 卷积核在此处的作用与对每个时间步使用一个全连接层（fully-connected layer）是等价的。  
数学上，对所有时间步  $t$  的操作可以写成：

$$\mathbf{M} = \{\mathbf{m}(1), \mathbf{m}(2), \dots, \mathbf{m}(T)\} \in \mathbb{R}^{2 \times T},$$

对于一个批次中的所有样本，则有输出张量为  $\mathbb{R}^{B \times 2 \times T}$ 。

2.2 激活函数与数值稳定性

在实际使用时，为了保证掩码值在  $[0, 1]$  范围内，通常会在卷积操作后接一个非线性激活函数（如 sigmoid），将线性输出映射到 (0,1) 内。例如：

$$\hat{\mathbf{m}}(t) = \sigma(\mathbf{W} \mathbf{f}(t) + \mathbf{b}),$$

其中  $\sigma(z) = \frac{1}{1+e^{-z}}$  是 sigmoid 函数。这样做可以确保掩码的物理意义，即每个掩码值表示对应频率或时刻上某个语音成分的“保留”程度。  
在模型代码中，虽然在这一层的定义中没有直接写出 sigmoid，但在前向传播中通常会对 `self.mask_predictor(separated_feat)` 的输出调用 `torch.sigmoid` 函数，确保输出满足所需范围。

3. 代码实现解析

代码实现部分非常简洁：

```
self.mask_predictor = nn.Conv1d(encoder_dim, 2, kernel_size=1)
```

这一行代码实现了以下几点：

- 1. **输入输出通道设置**
  - 输入通道数： `encoder_dim` 。这意味着该卷积层接收来自分离网络的特征，其每个时间步由 `encoder_dim` 维度描述。
  - 输出通道数： `2` 。代表预测出两个掩码，一个对应每个说话人。
- 2. **卷积核大小**
  - `kernel_size=1` 表示该卷积操作仅在通道维度上进行线性组合，而不涉及时间步的局部信息混合。数学上，这相当于对每个时间步独立地执行一次矩阵乘法  $\mathbf{W} \mathbf{f}(t) + \mathbf{b}$ 。
- 3. **无需空间（时间步）降采样**
  - 由于卷积核大小为 `1` 且没有设置步长或填充参数，输出的时间步数与输入保持一致，保证了每个时间步都获得了相应的掩码预测。
- 4. **后续处理**
  - 通常在前向传播时，会对该卷积层的输出再使用 `sigmoid` 函数（例如 `masks = torch.sigmoid(self.mask_predictor(separated_feat))` ），将线性输出映射到 `[0, 1]` 的区间内，确保掩码的值适合作为加权系数使用。

4. 掩码预测在语音分离中的作用

掩码预测模块在整个语音分离流程中扮演着“选择性放大”与“抑制”特征的关键角色。其具体作用包括：

- **信息加权**

预测出的掩码  $\hat{\mathbf{m}}(t)$  会与原始的 STFT 幅度谱（或时域编码特征）进行逐元素相乘，形式上可以写为：

$$\tilde{X}(t, f) = \hat{m}(t) \cdot X(t, f),$$

其中  $X(t, f)$  为频谱幅度， $\hat{m}(t)$  是在当前时间步的掩码值。通过这种乘法操作，高于 0.5 的掩码值会保留更多信号能量，而较低的掩码值则会抑制干扰噪声或非目标语音。

- **多说话人分离**

当输出通道为 `2` 时，模型可以为每个说话人分别生成一个掩码。假设  $\hat{m}_1(t)$  和  $\hat{m}_2(t)$  分别对应两个说话人的掩码，那么最终可以得到两个独立的加权频谱表示：

$$\tilde{X}_1(t, f) = \hat{m}_1(t) \cdot X(t, f), \quad \tilde{X}_2(t, f) = \hat{m}_2(t) \cdot X(t, f).$$

这种方法使得后续的逆变换（如 iSTFT）能够分别重构出两个不同的语音信号，从而实现语音分离的目标。

- **线性映射的优点**

使用 `kernel_size=1` 的卷积层作为掩码预测器的好处在于，计算简单、参数较少，而且能在不引入额外时序信息混合的情况下直接对每个时间步的特征向量进行全局线性变换。数学上，这样的操作保证了时间步之间的独立性，防止了因跨时步信息混合而导致的非必要干扰。

5. 数学优势与实际意义

从数学角度看，掩码预测模块将原始特征  $\mathbf{f}(t)$  映射到一个低维的掩码向量空间，通过线性变换加上非线性激活（如 `sigmoid`），实现了以下几点：

- 1. **降维与特征压缩**

输入特征  $\mathbf{f}(t) \in \mathbb{R}^C$  被映射到  $\mathbb{R}^2$  空间，降低了维度。这种降维过程不仅减少了模型参数，还使得后续信号重构过程更加高效和明确。
- 2. **局部独立性**

由于卷积核大小为 `1`，该映射对每个时间步是独立进行的。这保证了每个时间步的掩码预测只依赖于该时刻的特征向量，不受其它时刻信息的干扰，有利于捕捉瞬时的语音变化。
- 3. **线性与非线性结合**

线性变换部分使得模型能够捕捉到输入特征的全局线性关系，而后续通过 `sigmoid` 等非线性函数的映射，进一步保证了输出值在合理范围内，并引入必要的非线性特性，使得掩码具有更好的判别能力。
- 4. **简单高效**

从计算复杂度角度，该层只需进行一次  $2 \times encoder\_dim$  的矩阵乘法，相较于更复杂的网络结构，其计算开销极低，非常适合在实时或资源受限的环境中使用。

### 3. 前向传播（Forward Pass）

在本模型中，前向传播过程将输入原始波形经过多个模块依次处理，最终生成每个说话人的分离语音。下面我们从时域特征提取、频域转换、时频融合、掩码预测以及频域掩码应用与逆 STFT 还原五个部分对前向传播流程进行详细解析。

#### 3.1 时域特征提取

```
time_feat = F.relu(self.time_encoder(waveform))
```

##### 数学与代码解析

- 输入与卷积操作**  
输入的波形  $\mathbf{x} \in \mathbb{R}^{B \times 1 \times L}$ （其中  $B$  为批量大小， $L$  为采样点数）经过时域编码器（一个一维卷积层）处理。时域编码器的作用是利用一组卷积核（视为学习到的基函数）将原始信号映射到一个更高维且稀疏的潜在表示空间。设卷积核权重为  $\mathbf{W}_{\text{time}} \in \mathbb{R}^{\text{encoder\_dim} \times 1 \times k}$ （其中  $k$  为卷积核大小），对于每个输出时间步  $t$ ，卷积运算为：

$$\text{time\_feat}[b, c, t] = \max \left( 0, \sum_{i=0}^{k-1} W_{\text{time}}(c, 0, i) \cdot x(b, 0, t + i - P) + b_{\text{time}}(c) \right),$$

其中  $P$  为填充量，ReLU 激活函数 ( $\max(0, \cdot)$ ) 引入非线性，使得负值被截断。最终输出的特征张量  $\text{time\_feat} \in \mathbb{R}^{B \times \text{encoder\_dim} \times T_{\text{time}}}$  表示在时域上对局部波形的响应。

#### 3.2 频域转换与特征提取

```
spec = torch.stft(waveform.squeeze(1), n_fft=self.n_fft, hop_length=self.hop_length, return_complex=True)
spec_mag = spec.abs().unsqueeze(1) # 计算幅度谱
freq_feat = F.relu(self.freq_encoder(spec_mag.squeeze(1)))
```

##### 数学与代码解析

- STFT 变换**  
为了获取频域信息，对输入波形进行短时傅里叶变换（STFT）。先通过 `waveform.squeeze(1)` 将形状从  $\mathbb{R}^{B \times 1 \times L}$  转换为  $\mathbb{R}^{B \times L}$ ，然后计算 STFT 得到复数谱图：

$$X(b, k, t) = \sum_{n=0}^{n_{\text{fft}}-1} x(b, n + tH) w(n) e^{-j \frac{2\pi kn}{n_{\text{fft}}}},$$

其中  $k = 0, 1, \dots, \frac{n_{\text{fft}}}{2}$  表示频率 bin 数（共  $\frac{n_{\text{fft}}}{2} + 1$  个）， $t$  表示时间帧， $H$  为 `hop_length`。

- 幅度谱计算**  
对复数谱图取模 (`abs`)，获得幅度谱  $|X(b, k, t)|$ ；随后使用 `unsqueeze(1)` 在通道维度上扩展，得到形状为  $\mathbb{R}^{B \times 1 \times (\frac{n_{\text{fft}}}{2} + 1) \times T}$  的张量。
- 频域特征提取**  
使用 `squeeze(1)` 将通道维度去除，使得幅度谱变为  $\mathbb{R}^{B \times (\frac{n_{\text{fft}}}{2} + 1) \times T}$ 。接着，将该张量作为输入传入频域编码器（1D 卷积层），该层的卷积核将沿着时间轴对各频率 bin 进行线性组合。数学上，对于每个时间步  $t$  和每个输出通道  $c$ ，卷积操作为：

$$\text{freq\_feat}[b, c, t] = \max \left( 0, \sum_{k=0}^{K-1} W_{\text{freq}}(c, k) \cdot |X|(b, k, t) + b_{\text{freq}}(c) \right),$$

其中  $K = \frac{n_{\text{fft}}}{2} + 1$ （输入通道数），激活函数同样为 ReLU。输出的  $\text{freq\_feat} \in \mathbb{R}^{B \times \text{encoder\_dim} \times T}$ 。

### 3.3 时频特征融合

```
time_feat = F.interpolate(time_feat, size=freq_feat.shape[2], mode="nearest")
combined_feat = torch.cat([time_feat, freq_feat], dim=1)
fused_feat = self.domain_fusion(combined_feat)
```

#### 数学与代码解析

- **对齐时间步数**  
为确保时域特征 `time_feat` 与频域特征 `freq_feat` 在时间维度  $T$  上一致，使用插值函数 `F.interpolate` 将 `time_feat` 调整为与 `freq_feat` 相同的时间步数。记调整后的时域特征仍为  $\text{time\_feat} \in \mathbb{R}^{B \times \text{encoder\_dim} \times T}$ 。
- **特征拼接**  
将时域与频域特征在通道维度上拼接：

$$\text{combined\_feat} = \text{Concat}(\text{time\_feat}, \text{freq\_feat}) \in \mathbb{R}^{B \times (2 \times \text{encoder\_dim}) \times T}.$$

数学上，这相当于将每个时间步的特征向量从  $\mathbb{R}^{\text{encoder\_dim}}$  扩展到  $\mathbb{R}^{2 \times \text{encoder\_dim}}$ 。

- **融合操作**  
通过时频融合模块（一个由 1D 点卷积和 ReLU 激活构成的子网络）对拼接后的特征进行降维与重组。该模块实现的映射为：

$$\text{fused\_feat} = \phi(\text{combined\_feat}) \in \mathbb{R}^{B \times \text{encoder\_dim} \times T},$$

其中  $\phi$  表示连续两层 1D 卷积（`kernel_size=1`）与 ReLU 激活的复合非线性映射。这样，融合后的特征既包含时域局部波形信息，也包含频域能量分布信息。

### 3.4 预测掩码

```
separated_feat = self.separation_net(fused_feat)
masks = torch.sigmoid(self.mask_predictor(separated_feat))
```

#### 数学与代码解析

- **分离网络**  
将融合后的特征  $\text{fused\_feat} \in \mathbb{R}^{B \times \text{encoder\_dim} \times T}$  输入分离网络。分离网络由两层 1D 卷积（`kernel_size=3`, `padding=1`）和 ReLU 激活构成，其作用是进一步提取局部时序特征：

$$\text{separated\_feat} = \psi(\text{fused\_feat}) \in \mathbb{R}^{B \times \text{encoder\_dim} \times T},$$

其中  $\psi$  表示该网络实现的非线性映射。

- **掩码预测**  
随后，利用掩码预测模块（一个 `kernel_size` 为 1 的 1D 卷积层）对经过分离网络处理的特征进行线性变换：

$$\mathbf{m}(t) = \mathbf{W} \mathbf{f}(t) + \mathbf{b} \quad \text{for each } t,$$

将每个时间步  $t$  的  $\text{encoder\_dim}$  维特征向量  $\mathbf{f}(t)$  映射到一个 2 维向量，分别对应两个说话人的掩码。通过对线性输出应用 sigmoid 函数：

$$\hat{\mathbf{m}}(t) = \sigma(\mathbf{W} \mathbf{f}(t) + \mathbf{b}),$$

确保掩码值位于  $[0, 1]$  范围内。最终输出的掩码张量形状为  $\text{masks} \in \mathbb{R}^{B \times 2 \times T}$ 。

### 3.5 频域掩码应用 & iSTFT 还原

```
masked_spec = spec.unsqueeze(1) * masks.unsqueeze(2)
reconstructed = torch.istft(masked_spec.reshape(B * C, FR, T), n_fft=self.n_fft, hop_length=self.hop_length, length=waveform.shape[0])
separated_speech = reconstructed.reshape(B, C, -1)
```

#### 数学与代码解析

- 掩码与频谱相乘**  
原始 STFT 结果  $\text{spec} \in \mathbb{C}^{B \times FR \times T}$  (其中  $FR = \frac{n_{\text{fft}}}{2} + 1$  为频率 bin 数) 先通过 `unsqueeze(1)` 扩展出一个新的维度, 变为  $\mathbb{C}^{B \times 1 \times FR \times T}$ 。同时, 将预测的掩码  $\text{masks} \in \mathbb{R}^{B \times 2 \times T}$  通过 `unsqueeze(2)` 扩展为  $\mathbb{R}^{B \times 2 \times 1 \times T}$ 。两者逐元素相乘:

$$\text{masked\_spec}[b, c, k, t] = \text{spec}[b, k, t] \times \hat{m}_c(t),$$

- 得到形状为  $\mathbb{C}^{B \times 2 \times FR \times T}$  的掩码应用后的频谱, 其中  $c = 1, 2$  分别对应两个说话人。
- 逆 STFT (iSTFT) 还原**  
为了将频谱还原回时域信号, 需要对每个掩码应用后的频谱进行逆 STFT。首先, 将张量 `reshape` 成形状  $\mathbb{C}^{(B \times 2) \times FR \times T}$ , 这样可以对每个说话人独立地进行 iSTFT。逆 STFT 操作数学上为:

$$\tilde{x}(b, c, n) = \text{iSTFT} \left( \text{masked\_spec}_{(b,c)} \right),$$

- 其中  $n$  表示时域采样点, 参数  $n_{\text{fft}}$ 、`hop_length` 和 `length` (保证还原后的采样点数与原始波形相同) 保持一致。输出为时域信号, 形状通常为  $\mathbb{R}^{(B \times 2) \times L}$ 。
- 重塑输出**  
最后, 将重构后的信号 `reshape` 为  $\mathbb{R}^{B \times 2 \times L}$ , 其中每个批次中包含两个说话人的语音信号, 即:

$$\text{separated\_speech} \in \mathbb{R}^{B \times 2 \times L}.$$

### 整体流程总结

- 时域特征提取**  
使用一维卷积和 ReLU 从原始波形中提取时域特征, 得到 `time_feat`。
- 频域转换与特征提取**  
对原始信号进行 STFT 得到频谱, 计算其幅度谱后, 再通过一维卷积提取频域特征, 得到 `freq_feat`。
- 时频特征融合**  
对齐时域与频域特征的时间步数后, 将两者在通道维度上拼接, 并利用点卷积进行融合, 得到综合时频信息的融合特征 `fused_feat`。
- 预测掩码**  
将融合后的特征输入分离网络, 经过非线性映射后利用一个 1D 卷积层预测出两个说话人的掩码, 并通过 `sigmoid` 将掩码值限制在  $[0, 1]$ 。
- 频域掩码应用与还原**  
将预测出的掩码应用于原始 STFT 得到的频谱, 对每个说话人的频谱进行加权, 之后通过 iSTFT 将加权后的频谱还原为时域语音信号, 最终输出形状为  $\mathbb{R}^{B \times 2 \times L}$  的分离语音。