

# 语音增强初探

刘文哲



# Contents

## Part I 基于深度学习的语音增强

<b>1</b>	<b>基于幅度谱的语音增强 .....</b>	<b>9</b>
1.1	<b>基于掩蔽的语音增强算法 .....</b>	<b>12</b>
1.1.1	<b>处理流程 .....</b>	<b>12</b>
1.1.2	<b>用于幅度谱的时频掩蔽 .....</b>	<b>15</b>
1.2	<b>基于谱映射的语音增强算法 .....</b>	<b>17</b>
1.2.1	<b>处理流程 .....</b>	<b>17</b>
1.3	<b>小结 .....</b>	<b>22</b>
	<b>References .....</b>	<b>22</b>
<b>2</b>	<b>基于复数谱的语音增强 .....</b>	<b>23</b>
2.1	<b>基于复数谱映射的语音增强 .....</b>	<b>24</b>
2.2	<b>基于笛卡尔坐标系的复值掩蔽方法 .....</b>	<b>27</b>
2.3	<b>基于极坐标系的复值掩蔽方法 .....</b>	<b>31</b>
2.4	<b>基于幅度相位分量估计的语音增强 .....</b>	<b>33</b>
2.5	<b>小结 .....</b>	<b>35</b>
<b>3</b>	<b>混合式和生成式语音增强模型 .....</b>	<b>37</b>
3.1	<b>混合式语音增强模型 .....</b>	<b>37</b>
3.2	<b>生成式语音增强模型 .....</b>	<b>40</b>
3.3	<b>小结 .....</b>	<b>40</b>
<b>4</b>	<b>基于解耦的语音增强 .....</b>	<b>41</b>
<b>5</b>	<b>基于时域的语音增强 .....</b>	<b>43</b>
<b>6</b>	<b>语音增强算法的因果性 .....</b>	<b>45</b>
6.1	<b>语音音量的归一化 .....</b>	<b>46</b>
6.2	<b>因果模块讨论 .....</b>	<b>47</b>



## **Part I**

### 基于深度学习的语音增强

语音增强算法引入深度学习大抵是一种必然，讨论基于深度学习的语音增强算法我们必须首先把目光放在Prof. Wang的OSU-PNL身上。Prof. Wang 是用于语音分离的计算听觉场景分析(CASA)领域的代表人物，他提出用理想二值掩蔽(IBM)解决语音增强(分离)问题。理想二值掩蔽，其中掩蔽广泛用于图像处理之中，具体到语音增强领域是指对选定的时频区域进行遮挡以控制要处理的区域，二值则是对每个时频点量化的精度非0即1，理想二值掩蔽的含义则是根据纯净语音和噪声之间的能量关系，将语音能量占主导的时频点标记为1，噪声能量占主导的时频点标记为0，由此得到一个滤波器，对带噪特征进行滤波，在语音时频稀疏性假设下留下的即使语音成分。于是，熟悉机器学习的读者应该很容易发现，语音增强问题在IBM的定义下可以被看做一个分类问题，当时OSU-PNL的Y. Wang已经开展了支持向量机(SVM)估计IBM的研究。机器学习中特征提取和目标选择十分重要，作为机器学习模型的SVM当然也不例外。Y. Wang的研究当然也包括评估众多利用听觉特性构造的手工特征对SVM进行语音增强的性能影响，在深度学习大行其道的今天，利用深度学习提取更有效的特征的思维并不难理解，而Y. Wang在2013年也的确选择这样去做。如果你对机器学习有一定了解的话，应该了解机器学习任务可以分为分类任务和回归任务两种，基于回归模型的语音增强算法也于2014年被中科大USTC SPRAT的Y. Xu等人提出。基于回归的语音增强模型直接利用神经网络建立从带噪语音谱特征到纯净语音谱特征的映射，利用网络直接生成增强的对数功率谱。如果说Y. Wang的工作是基于监督学习的CASA算法的延伸的话，Y. Xu的这项工作则更具有基于统计模型的传统语音增强算法的风格。工作中沿用了MMSE-LSA中的对数谱特征以及传统语音增强方法中人耳对相位信息不敏感的假设。然而利用分类和回归去分类这两个极具代表性的工作并不利于对后续方法的研究，因为后续的理想比值掩蔽(IRM)、相位敏感掩蔽(PSM)等训练目标也明显受到了传统语音增强统计模型的影响，这类掩蔽估计由于其定义也不再能被视为分类问题(而同样是回归问题)。这类掩蔽本质上是在设计某种意义上最优的滤波器，毕竟，回顾IBM的介绍“IBM的含义正是根据每个时频点上语音和噪声之间的能量关系，将语音能量占主导的时频点标记为1，噪声能量占主导的时频点标记为0”中“主导”二字充满了不安全感。一个极端的例子是，在极低信噪比情况下IBM会出现非常多的0值导致理想掩蔽的滤波结果的语音质量和可懂度都会十分糟糕。IRM将IBM的硬分类变成了一种软分类，每个时频点的滤波器系数定义为语音能量和语音噪声二者能量之和的比值。熟悉传统语音增强算法的读者不难发现，这个定义和频域维纳滤波算法尽管并不相同但还是十分相像。因此，Prof. Wang在综述中从训练目标角度将上述两大类工作分别定义为基于掩蔽(masking-based)和基于谱映射(mapping-based)的算法用于分类，这种分类方式得到了广泛地接受。

尽管上述内容足以概括基于深度学习的语音增强在早期阶段的历史，然而除了训练目标，模型架构的发展同样值得我们关注。不过在这一段中我们对文献中的训练目标使用掩蔽还是谱特征并不感兴趣，因此，如无必要在介绍中不会涉及映射目标。另一方面，虽然本节介绍的是深度学习语音增强的发展，然而为了保证叙述的连贯性，本段刻意忽略了算法的提出时间，只保证了相关联工作的顺序。由于语音具有时序性，因此时间建模是网络结构设计上不可忽视的问题。前面提及的方法都是利用全连接网络，考虑时序关系的方式也十分朴素——将当前帧前后几帧组合在一起形成上下文窗(context

window)后送入网络。很明显，这种方式并不足以建模长时信息。而循环神经网络(RNN)天然适合时序建模，因此2017年OSU-PNL的J. Chen和USTC-SPRAT的Sun分别开展了基于LSTM的语音增强模型的工作。这类网络和全连接网络唯一结构上的区别在于全连接网络的其中几层(通常是前几层)被替换成LSTM层，但其在增强性能、未见说话人和噪声类型的泛化能力上都明显优于全连接网络。RNNs的提出使得目前基于深度学习的语音增强模型几乎全部放弃了全连接网络。另一个需要提及的内容在于卷积神经网络的应用，熟悉图像分割的读者应该已经感受到了语音增强与图像分割的相似之处，因此卷积神经网络的引入也十分自然。不过在介绍卷积神经网络在语音增强的应用之前，有必要至少提一点语音增强与图像分割的区别，尽管我们可以将语音特征看成一张图片，但是由于语音增强有很多实时应用场景，因此在这张“图片”时间维度的处理需要十分谨慎“实时陷阱”。卷积层由于“局部连接、权重共享”的特点相比全连接层和循环层拥有更少的参数，2015年清华SATLab的L. Hui提出的maxout卷积网络首次在语音增强领域性能超过全连接网络。而后中央研究院Bio-ASP Lab的S. Fu和卡内基梅隆大学的S. R. Park分别将图像领域常用的普通卷积神经网络和卷积编解码器结构成功应用到了语音增强中。由于卷积编解码器结构对于时频点级预测任务的天然适配，基于卷积编解码器的语音增强模型逐渐成为主流。自然地，将卷积层的特征提取能力和循环层的时序建模能力结合的工作应运而生。2018年H. Zhao在微软亚研完成了卷积层特征提取、LSTM时序建模、最后由全连接层进行幅度谱估计的EHNet，同年OSU-PNL的K. Tan将LSTM插入卷积编解码器之间提出一种卷积循环网络(CRN)。除了RNN适合序列建模外，时序卷积网络(TCN)利用一维空洞卷积实现了并行时序建模过程。2018年K. Tan引入门控(gating)TCN取代RNN构造了一种全卷积增强网络。

随着基于深度学习的幅度谱估计算法性能的快速发展，两个很有意思的现象出现了。一是相对简短可以介绍完的，是深度学习反求诸传统，深度学习作为一种参数估计器为传统语音增强中的MMSE、MMSE-LSA和卡尔曼滤波等算法提供依赖先验假设更少且更为准确参数估计。一般将这类深度学习与传统算法结合的方法称为混合式模型，然而混合式模型类工作相比性能提升更重要的贡献在于让基于深度学习的语音增强渐渐褪去纯机器学习或深度学习意味的工作，而逐渐带有语音信号处理的特色。当然，相比于上述这种“主动”结合语音特色，另一种相对“被动”的结合方式产生于另一个有意思的现象——基于深度学习语音增强的性能第一次遇到了瓶颈。当然人们很快意识到了问题——这一性能限制是只进行幅度谱估计带来的，相位恢复似乎对人耳感知是重要的。以至于解决该问题的方案出现得如此之快在今天看来似乎并没有瓶颈的意味在。在相位恢复重要性思潮下，用于同时进行幅度估计和相位恢复的方案大方向上可以分为时域的方案和谱域的方案。不论是哪种方案他们都是为了解决由于相位谱没有清晰的模式结构导致其不易优化的问题。首先介绍时域模型，原因依旧是因为介绍起来可以比较简略。时域模型的兴起主要由于其在语音分离领域大获成功，以一种波形到波形的映射方式提供了一种解决相位谱不易优化问题的方案——放弃建模相位。不过也许是由于噪声模式的复杂性，抑或是房间冲激响应的敏感多变，时域模型在噪声混响场景下的性能和鲁棒性使其在语音增强中并不能成为主流。网络结构上和之前的幅度谱架构基本还是类似的，有卷积编解码器结构也有卷积循环(或TCN)网络，只不过对应的卷积层大多用一维卷积代替了二

维卷积(当然也同样有使用二维卷积作为编解码器的架构)。相比之下，谱域的方案解决相位估计的方式则略显复杂。众所周知，复数可以由幅度相位表示也可以由实部虚部表示。其中实部虚部表示和之前幅度谱模型一样，主要分为基于掩蔽和基于谱映射的方式。2016年OHU-PNL的Williamson提出了一种复数比值掩蔽(cIRM)，分别估计实部和虚部的掩蔽；而2017年S. Fu提出了一种谱映射的CNN，2019年K. Tan将CRN应用到复数谱映射直接利用网络估计增强语音的实部和虚部。2019年首尔国立大学的H. S. Choi提出了一种基于极坐标系的复值掩蔽策略(pcRM)，通过网络估计的cRM得到有界的幅度掩蔽和无界的相位掩蔽。尽管这篇文章成功应用了复值网络并且西工大ASLP的Y. Hu继复值U-Net发展为复值CRN取得了令人瞩目的性能，然而本篇文章对于在不同损失函数情况下的几种典型复值掩蔽的分析可能更加有趣：尽管这类复数目标算法都宣称致力于相位的恢复，然而其带来的性能提升究竟是源于相位的恢复还是噪声成分幅值的抑制(相位可能调整的极少)值得探究，D. Yin等人同样通过实验发现了使用cRM虚部几乎为0的问题。而为了体现相位的修复，基于幅度相位表示的方案得到了考虑，其中包括采用双分支结构同时估计幅度和相位，通过多任务学习的方式达到相位恢复的方法；也有一些算法采用了更具结构性的相位目标与幅度谱同时和先后估计，比如OHU-PNL的Z. Wang提出的群时延(GD)以及西工大CIAIC的N. Zheng提出的瞬时频率倒数(IFD)。

尽管Z. Wang的分析要晚于接下来要介绍的一些工作，然而在此处提及也许是必要。尽管这篇分析的重点在于解释RI+Mag损失函数为何是有效的，但其提供的幅度估计和相位估计的补偿作用无疑为解耦风格工作有效性提供了一种合理解释。所谓幅度和相位估计的补偿作用，是指在迫使网络输出逼近理想复数谱时，幅度谱估计的准确度会大打折扣。而如Griffin-Lim相位重构的很多相位估计算法都在幅度谱估计的基础上进行。一个解决幅度相位估计补偿作用的想法是，将复数谱的优化问题解耦为幅度初估计和包括相位信息在内的“全信息”估计两部分。尽管这仍属于复数谱估计的一类方法，然而这类算法在近年来备受关注且性能优异，因此有必要单独作为成段介绍。2020年，奥尔登堡大学通信声学的N. L. Westhausen先后利用两次信号变换——首先利用STFT幅度谱上进行幅度估计，而后再利用一维卷积生成的可学习的分析/合成基函数直接将语音时域波形映射到高维空间进行语音和噪声成分的分离，利用时域模型的方式完成相位恢复。哈工大SPLab的Z. Du在Mel谱(一种符合人耳听觉特性的幅度谱变换)域进行降噪，而后利用语音合成声码器根据增强的Mel谱合成语音波形从而实现相位信息的引入。中科院声学所IACAS-Lab9的A. Li则在2021年提出一种完全在STFT域上先后完成幅度谱估计和复数谱残差修正的框架，利用复数谱残差修复相位信息。完全在STFT域上处理两部分任务的一个好处是，原本级联的两个子任务出现了并行优化的可能，2022年A. Li和西工大ASLP的Y. Fu分别提出了幅度估计和复数谱修正的两种并行架构。其实，从复数谱的角度考虑，这种解耦是由于幅度相位估计的补偿作用导致幅度谱估计不准确造成的，因此，除了首先约束幅度谱后再根据相位差微调增强结果，也可以先得到复数谱结果后补偿不准确的幅度信息。于是在2022年，北交大IIS的T. Wang则先进行复数谱估计再利用谐波特征进行幅度谱修正的方法。而解耦风格的工作也并非只由从训练目标解耦一种方式。从语音被噪声污染的过程角度，早在2016年USTC-SPRAT的T. Gao就提出以不同信噪比的被污染语音作为目标将一次降噪过程

分解为渐进地、多阶段地、逐信噪比地提升。从语谱高低频特性差异角度，搜狗的J. Li对语谱进行子带分解，多阶段地处理各个子带。而对语音特性进行分析，语音可以分解为包络和周期两部分，2018年，J.-M. Valin首先利用GRU估计Bark域(另一种符合人耳听觉特性的幅度谱变换)的增益因子，而后利用信号处理的方式进行基频滤波抑制谐波间噪声残留，而后在2020年，其又对模型中各部分进行进一步提升，以可观的复杂度达到超过众多复数谱模型的性能。而对该工作神经网络化最杰出的代表当属PAU模式识别实验室的H. Schröter提出的在估计包络增益因子外利用网络估计深度滤波器(deep filtering)系数的框架。除此之外，语音增强任务也可以理解为噪声去除和语音恢复两部分，因为单通道语音增强具有噪声抑制和语音失真之间的折衷问题，彻底去除噪声势必会对语音造成严重损伤。因此，一种先进行激进地噪声抑制而后再恢复在降噪过程中的失真语音的框架首先由TU Braunschweig通信技术所信号与机器学习组的M. Strake于2019年提出，2020年内蒙古大学IMUAI的X. Hao在该思路的基础上引入计算机视觉的概念形象地将该过程命名为“masking and inpainting”即将所有噪声的时频点去除后根据语谱相关性修补被破坏的语音时频点。而传统的谱减法则从带噪语音的合成过程反推，将语音增强看做噪声估计和移除噪声成分的过程。2020年，哥伦比亚大学C2G2的R. Xu受传统谱减法的启发将降噪过程分解为噪声的“masking and inpainting”过程(利用VAD得到静音段而后通过静音段推测出完整的噪声谱)并将噪声信息和带噪语音信息一同作为输入完成降噪。2021年，IACAS-Lab9的W. Liu将基于训练目标解耦的思路与基于带噪语音合成的解耦思路结合首先利用多任务学习幅度域对语音降噪并直接估计噪声成分，而后修正复数谱恢复过度抑制的语音。需要注意的是，尽管本段从解耦角度介绍了以上算法，但是这些算法在提出时受到的启发性工作和动机可能并非如此。这种分类目前也并不作为一种共识，而更多的是一家之言。

需要注意的是，网络架构和学习目标虽然的研究虽然更多然而只是性能提升的一方面，更长时间范围的输入特征以及精心设计的损失函数对模型性能的提升可能更加有益，其中损失函数由于不会在推理过程中引入额外的时延和复杂度而备受关注。一个有力的证据是在ICASSP2021 DNS-Challenge中微软提供的结合一种幂律压缩的复数谱均方误差损失的复数掩蔽网络相比包括解耦风格框架在一众高参数量和计算量方案仍表现出最优的性能。损失函数包括基于距离的损失函数、基于能量比的损失函数和基于感知指标的损失函数等。常见的距离表征最常用的莫过于基于范数的损失函数 $L_p$ 损失，即平均绝对误差(MAE)和均方误差(MSE)，这两类误差都被用于度量时域波形、掩蔽、频谱以及嵌入式特征。时域波形的 $L_p$ 度量形式相对简单，表示为 $\mathcal{L}_{time-L_p} = \|s - \hat{s}\|_p$ ，有时会加入对噪声项 $L_p$ 损失作为正则，以提升语音估计性能，有文章表示，对于时域样本而言MAE相比MSE更关注小信号并具有更强的噪声抑制能力。除了IBM有使用交叉熵损失进行度量外，其余掩蔽大多都会采用 $L_p$ 损失直接计算掩蔽间的误差 $\mathcal{L}_{mask} = \|M - \hat{M}\|_p$ ，J. Valin提出了幂律压缩的掩蔽MSE损失 $\mathcal{L}_{RNNoise} = \|M^c - \hat{M}^c\|_2$ ，利用指数项 $c$ 控制噪声抑制的激进程度。然而，目前普遍认为基于信号近似(SA)技术的损失相比直接度量掩蔽间差异对语音增强的性能更有益： $\mathcal{L}_{SA} = \|S - \hat{M} \cdot Y\|_p$ 。可以看出，此时对掩蔽的范数损失已经和基于频谱的范数损失基本一致，只需将滤波项 $\hat{M} \cdot Y$ 改为网

络映射的谱 $\widehat{S}$ :  $\mathcal{L}_{spec-L_p} = \|S - \widehat{S}\|_p$ 。当然 $S$ 和 $\widehat{S}$ 可以是幅度谱或其变换(如对数谱、Mel谱等), 也可以是复数谱。由于幅度谱和复数谱是如此常用, 这里我们分别给出其具体形式:  $\mathcal{L}_{Mag} = \mathbb{E}[|A - \widehat{A}|^p]$ 和 $\mathcal{L}_{com} = \mathbb{E}[|\mathcal{R}(X) - X(\widehat{X})|^p + |\mathcal{I}(X) - \mathcal{I}(\widehat{X})|^p]$ (其中 $A$ 和 $X$ 分别代表幅度谱和复谱)。 $\mathcal{L}_{spec-L_p}$ 还可能根据如AMR编码等心理声学机制进入加权, 以期望其更符合人耳听觉, 而复数谱范数损失又常与幅度损失结合以缓解幅度相位估计的补偿作用, 该损失被称为RI+Mag损失:  $\mathcal{L}_{RI+Mag} = \alpha\mathcal{L}_{Mag} + \beta\mathcal{L}_{com}$ 。在此基础上, 与时域波形 $L_p$ 损失结合的 $\mathcal{L}_{time-spec-L_p} = \mathcal{L}_{RI+Mag} + \mathcal{L}_{time-L_p}$ 、基于幂律压缩的RI+Mag损失 $\mathcal{L}_{cprs} = \alpha\mathbb{E}[|A^c - \widehat{A}^c|^2] + \beta\mathbb{E}[|\frac{X \cdot A^c}{A} - \frac{\widehat{X} \cdot \widehat{A}^c}{\widehat{A}}|^2]$ 以及多分辨率(假设有 $I$ 种STFT点数对应的 $I$ 个分辨率)的 $\mathcal{L}_{cprs}$ 函数 $\mathcal{L}_{MR-STFT} = \sum_i \mathcal{L}_{cprs}^i, i = 1, \dots, I$ 都成为目前广泛使用的 $L_p$ 损失。另外由于 $\mathcal{L}_{cprs}$ 对噪声的抑制能力较强导致语音失真较多, 一种非对称损失有时也被作为正则项 $\mathcal{L}_{asym} = \mathbb{E}[|\max(A^c - \widehat{A}^c, 0)|^2]$ 。可以看出, 尽管有文献提及由于STFT频点分布更类似拉普拉斯分布因此建议使用 $\mathcal{L}_1$ 范数损失, 然而大多数工作更喜欢选择 $\mathcal{L}_2$ 范数损失。 $L_p$ 范数损失最后也见于对嵌入式特征的距离度量, 利用wav2vec、PANN等预训练模型生成嵌入式特征之间的距离作为谱距离损失的正则项。此外, 一些除了 $\mathcal{L}_p$ 范数作为距离度量, KL散度也有一些工作将其与 $\mathcal{L}_p$ 范数联合度量谱距离:  $\mathcal{L}_{KL} = \mathbb{E}[\widehat{X} \cdot \log(\frac{\widehat{X}}{X})]$ 。基于信号能量比的损失函数主要用于时域波形, 常见的有SDR(及其Log形式)和SI-SDR损失, 分别定义为 $\mathcal{L}_{SDR} = -10\log(\frac{|s|^2}{|\widehat{s}-s|^2})$ 和 $\mathcal{L}_{SI-SDR} = -10\log(\frac{|\xi \cdot s|^2}{|s-\xi \cdot \widehat{s}|^2})(\xi = \frac{s^T \widehat{s}}{s^T s})$ 。尽管SI-SDR被广泛应用于语音分离当中, 但是由于其过多的语音失真导致最近的工作常将其与 $L_p$ 范数谱距离损失联合使用。上述特征常被诟病并不能反映人耳的听感, 因此, 一些感知指标, 如PESQ、ESTOI和CD, 一度被引入作为损失函数训练模型。然而这类损失函数仅能提升被优化指标的性能, 并没有带来其余指标的提升。要知道, 目前的客观指标与人类的主观听感并不完全一致, 单纯提升某些客观指标大多也并未带来主观听感的明显提升。于是近年来, 这类感知指标类损失更多是以一种联合形式作为验证集损失被用于选择最优的模型。

尽管生成对抗网络(GAN)常作为一种不同的网络架构, 然而这里我们更希望将其作为一种训练手段, 即对抗性训练。众所周知, GAN由生成器和鉴别器两部分, 在GAN在语音增强的使用中, 其生成器可以用前文提到的众多网络代替, 而之前提到的各种损失同样可以作为生成器损失函数的其中一项或等价于生成器损失函数。因此, GAN在语音增强中目前表现的作用, 更类似于一种类损失函数”用于帮助网络专注于那些原本损失函数难以强调的噪声成分和语音伪影(artifacts)部分。于是, 一言以蔽之GAN在语音增强的整体发展: 之前的语音增强模型(如LSTM、UNet等)代入生成器, 各式GAN的发展技术(如LSGAN、Wasserstein GAN、Relativistic GAN、Conditional GAN和Geometric GAN等)代入鉴别器。如此导致GAN始终并未在语音增强领域至少在实时语音增强领域成为主流技术, 这很可能与语音增强任务与图像生成、语音合成等生成类任务的差异过大有关, 然而, 在无监督或自监督情况下, 包括GAN在内的生成式模型(还包括变分自动编码VAE等)也许是值得研究的方法。

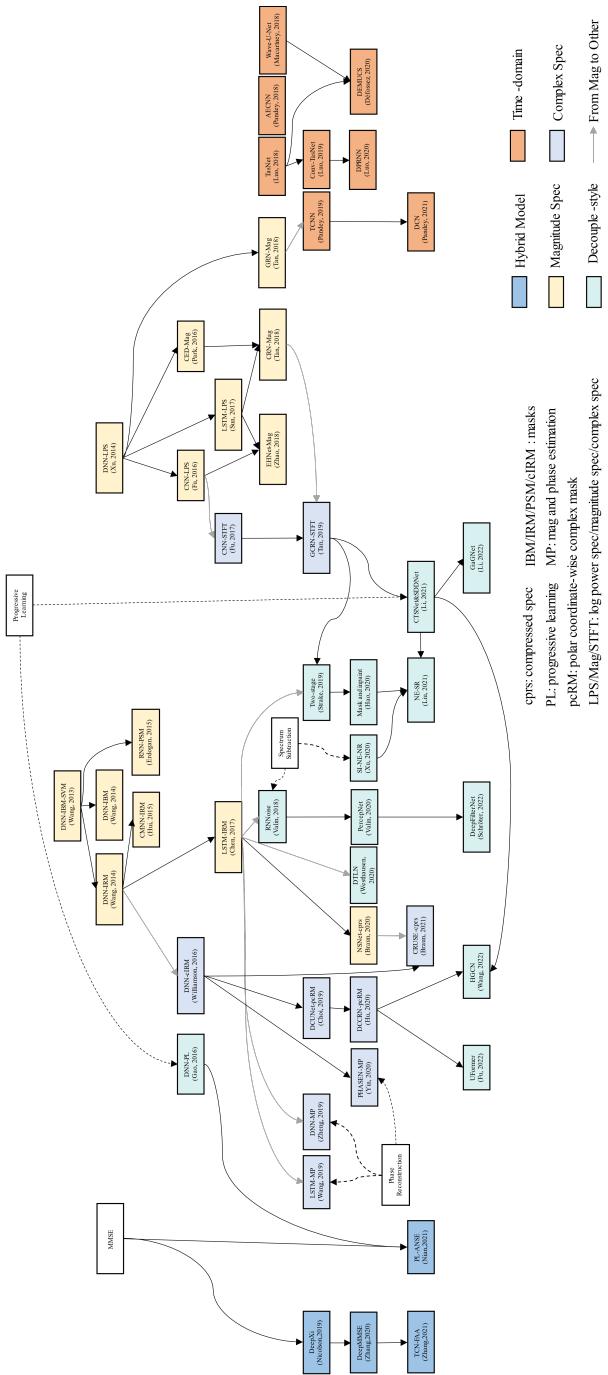


Fig. 0.1 深度学习语音增强算法发展

图0.1梳理了深度学习语音增强发展的部分脉络，然而对于目前处于发展之初的无监督语音增强、超宽带语音增强、多模态语音增强、个性化语音增强等并未涉及，去混响等问题也未讨论。

# Chapter 1

## 基于幅度谱的语音增强

由于人耳对于幅度相比比相位更加敏感，传统的语音增强算法大部分都集中于在幅度谱上进行处理。于是，深度学习最早被引入语音增强时也在幅度谱上最先应用。基于幅度谱的语音增强过程大抵可以按照图1.1所示的流程。带噪信号通过傅里叶变换基或基于人耳听觉的滤波器组将时域信号分解成二维频域表征，而后提取其幅度域特征，如短时傅里叶变换的幅度谱、对数功率谱等。通过纯净语音和噪声可以计算得到网络的映射目标，常见的是时频掩蔽或与特征提取对应的谱特征。特征提取和分离目标组成输入输出对用于对网络模型进行训练，建立从输入特征到分离目标之间的映射。训练好的模型作为分离模型在测试阶段通过接收混合语音时频分解后提取的特征预测出估计的分离目标。分离目标结合混合信号通过时频分解的逆变换即可得到增强的信号波形。

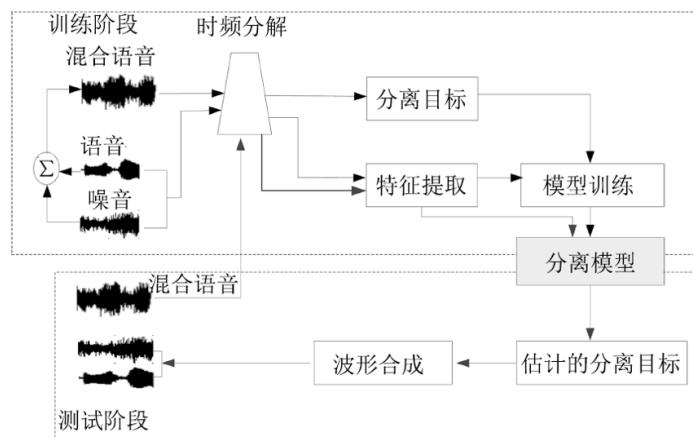
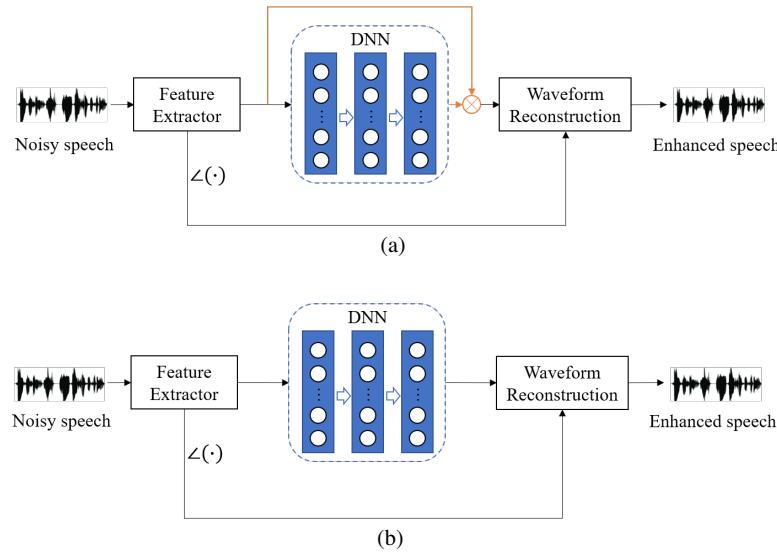


Fig. 1.1 基于幅度谱的语音增强框图[1]

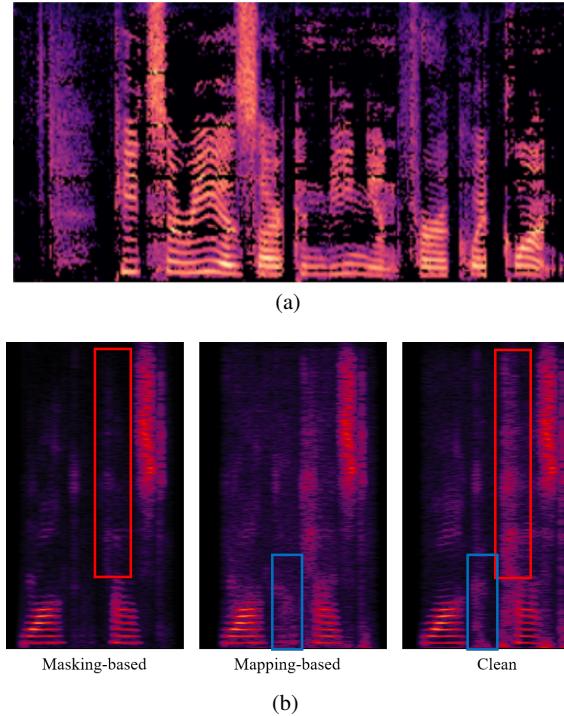
目前通常将基于幅度谱的语音增强方法分为基于掩蔽(masking-based)和基于谱映射(mapping-based)两大类。基于掩蔽的语音增强方法利用网络预测出一个滤波器对幅度谱进行滤波，而后将滤波后的幅度谱与带噪相位耦合并重构回时域波形；基于谱映射的语音增强方法利用神经网络直接建立从带噪谱特征到纯净谱特征之间的映射，而后将谱特征与带噪相位耦合得到增强的语音。这两类方法的框图如图1.2所示，可以看出，二者的区别主要在于是否有显式的滤波过程(橙色部分)。



**Fig. 1.2** (a) 基于掩蔽的语音增强处理流程，(b) 基于谱映射的语音增强处理过程。

由于基于掩蔽方法是估计滤波器并对带噪语谱滤波，早期这类算法估计的滤波器都是起抑制作用，不可避免地在语音能量弱但存在强噪声的谱上形成不连续的黑洞(black holes)，由于语音能量弱的地方可能是谐波之间(图1.3(a))、中高频谐波、以及轻音(图1.3(b))，这种不连续极其影响语音质量和可懂度；而基于谱映射的方法可以称之为脑补，算法的作用可以看成生成或者再生，凡是噪声存在的时频点均有可能完全抹去后重生出新的、结构类似的频谱，不过也正是由于谱映射的生成机制，早期的伪影(artifacts)问题(图1.3(b))、过平滑(over-smoothing)问题和鲁棒性问题常令人诟病。不过随着两类算法后来的发展，上述问题已经不再明显。

我们将图1.3(a)或(b)展示的语音的时频关系图称为语谱，语谱中的每个点被称为时频点(time-frequency bin)。



**Fig. 1.3** (a) 频谱黑洞[4], (b) 利用基于掩蔽和基于谱映射方法处理的语音, 基于掩蔽的方法的清音被吞掉(见红框), 基于谱映射的方法在无谐波部分脑补出谐波(见蓝框)。

需要注意的是, 这种最常用的分类方式是以训练目标作为分类依据的。由于语音增强的应用角度可分为面向人耳听感的和面向机器感知的。前者的目的在于提高人耳在通信过程中的语音质量和可懂度, 不舒适的噪声成分的抑制是其处理的重点, 而不易感知的语音失真是能过够容忍的; 而后者往往对语音失真更加敏感, 而可以容忍一定的噪声残留。因此, 训练目标和损失函数相比特征更为重要。很可能基于此这种分类方式成为主流。基于幅度谱的语音增强方法常采用两类特征: 幅度谱及其听觉感知的相关变换。前者直接采用短时傅里叶变换的幅度谱或者其对数/幂律压缩形式; 后者则将幅度谱变换到ERB/Bark/Mel/Gammatone等听觉感知, 有时也引入线性预测/基音预测/倒谱变换等额外处理得到一组听觉感知相关系数作为输入特征。随着深度学习模型的快速发展, 神经网络通过部分非线性变换取代了手工特征的提取过程, 使得对后者的研究日趋减少。因此, 我们不再对特征方面展开介绍, 如果读者希望深入了解, 可以参看[2]。

基于幅度谱的语音增强模型的发展于2013年由Wang等人首次提出, 他们将全连接网络应用于语音增强领域作为一个子带分类器估计IBM。2014年, Xu等人证明了全连接网络能够通过直接建立带噪语音到纯净语音之间对数功率谱的映射实现语音增强。需要注意的是, 这一时期的工作由于受到早期深度学习发展的限制保留了明显的时代特色, 诸如使用受限玻尔兹曼机

预训练初始化网络权重、对输入特征进行时序平滑滤、对大量手工提取特征的研究投入过多、以及一部分应对噪声、信噪比和说话人泛化问题的技术等，这些操作逐渐被新兴的深度学习技术取代而不再被使用。同时，为了让网络感知语谱的时序信息，这一时期的输入特征常将当前帧左右多帧(context window, 上下文窗)拼接后送入网络，这极大地增加了输入特征的维度，又不能很好地表征语谱的时序关系。因此，天然具有时序建模能力的RNN和LSTM被引入到面向语音识别的语音增强中被证明循环网络的有效性，Chen等人则将LSTM应用在面向听感的语音增强任务中并展示了LSTM对提升未见噪声和说话人泛化能力的作用，64通道的耳蜗谱被送入四层LSTM后预测得到IRM，从网络结构上改进了原本全连接网络的泛化问题(由于历史惯性，LSTM乃至CNN网络在早期也仍有很多工作输入上下文窗作为特征)。在1.1中，基于时频掩蔽的幅度谱语音增强方法将以GRU网络为例进行介绍。此外，同时估计时频掩蔽和谱信息后加权融合的多目标训练同样是这一时期的常见工作以期望结合两种训练目标的优点或应对平稳/瞬态等不同特性的噪声。通过图1.3可以发现，幅度谱具有清晰的时频模式，全连接和LSTM都没有很好的建模这种时频相关性。因此，Hui等人引入CNN估计幅度谱的IRM，在此基础上Fu等人利用多任务学习利用卷积神经网络同时估计纯净对数功率谱和SNR从而使网络隐式地关注SNR失配问题，卷积神经网络由三个卷积层和两层全连接层组成，该工作的一个重要贡献是实验证明了图像领域常用的池化层由于会丢失细节并不利于语音增强任务，目前语音增强中的卷积层已普遍不再级联池化层。一个自然的想法是结合CNN对时频相关性的建模和RNN的时序建模能力，于是卷积递归网络(CRN)被提出。常见的两种组合方式分别是卷积层组成的编码器和递归层组成的时序建模模块级联后再经过全连接层预测时频掩蔽或谱特征(如2018年的EHNet)，和在由卷积层和分卷积层分别组成的对称的卷积编解码器(Convolutional Encoder-Decoder, CED)之间插入递归层组成的时序建模模块的编码器-时序建模-解码器架构(如2018年的CRN)。后者对后续工作的影响更大，因此在1.2中选择了一个基于CRN的谱映射语音增强算法进行介绍。除了上述工作之外，应当注意的是这一时期多任务学习的方法已经被引入，在估计时频掩蔽或谱特性的同时，信噪比(Signal-to-Noise Ratio, SNR)、基频 $f_0$ 、语音活动检测(Voice Activity Detection, VAD)等与语音增强/分离任务相关的参数也会有选择地同时被估计以控制模型的泛化性能或用于后处理提升增强语音的质量。

## 1.1 基于掩蔽的语音增强算法

### 1.1.1 处理流程

为了更好的解释基于幅度谱的语音增强算法，我们选用NSNet作为基于掩蔽的语音增强算法的例子进行说明。

NSNet的框图如图1.4所示，对于初窥门径者也许这个框图有些复杂，不过我也并不打算阐述所有部分。请将图1.2(a)和图1.4对应起来，图1.4虚线以上的部分，从左至右的， $x(t)$ 代表图1.2(a)中的带噪波形(noisy speech)，从STFT至代表网络的"3-layer GRU (257) FC + Sigmoid (257)"之间的那部分是

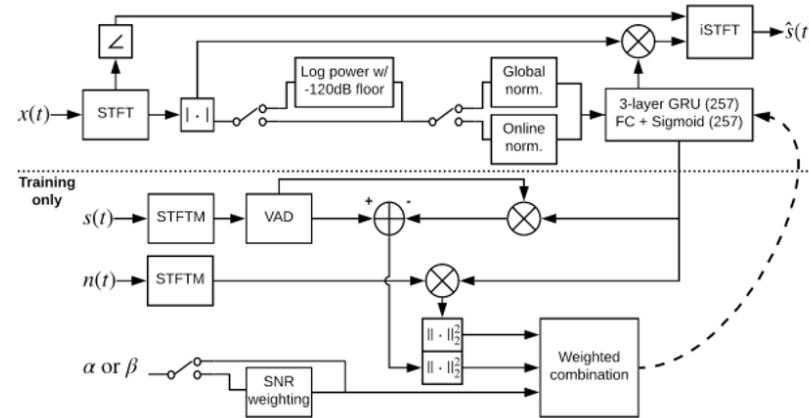


Fig. 1.4 NSNet框图[3]

特征提取器(feature extractor)，注意，图1.2(a)中的特征提取器有两个输出，一个送入网络(DNN)，另一个代表相位分量用于波形重构。请观察图1.4，同样可以发现相位从特征提取器中的STFT模块获得并传送给iSTFT模块，理所当然的，iSTFT对应的就是图1.2(a)中的波形重构(waveform reconstruction)。iSTFT的另一个输入自然是被网络估计的掩蔽滤波的特征。而 $\hat{s}(t)$ 则对应了图1.2(a)中的增强语音(enhanched speech)。至于虚线以下的部分代表的是损失函数的计算过程，与虚线以上的联系是图中唯一的虚线，用于表示通过损失函数计算得到的梯度的反向传播，以更新网络参数。总结一下，基于掩蔽的语音增强方法首先将带噪语音经过特征提取得到幅度谱相关的特征以及相位，而后幅度谱相关特征经过网络得到时频掩蔽，时频掩蔽通过滤波幅度谱(或幅度谱相关特征)得到增强的幅度谱(或幅度谱相关特征)。而后增强的幅度与带噪语音的相位结合重构得到增强的语音信号。NSNet采用了标准化的对数功率谱作为特征，而网络输出的是幅度谱的时频掩蔽。

输入特征不仅可以是对数功率谱、幅度谱、Mel谱和Bark域谱等都曾作为输入特征。输入特征的选取或多或少源于提出者的出身背景，比如对数谱、Mel谱和Bark域谱分别是传统语音增强MMSE-LSA、语音识别和音频编码领域常用的特征。

```

1 # 特征提取部分 wav_dataset.py lines 51-109
2 x_stft = torch.stft(noisy_waveform.view(-1), ...)
3 x_ps = x_stft.pow(2).sum(-1)
4 x_lps = LogTransform()(x_ps)
5 x_ms = x_ps.sqrt()
6 ...
7 for frame_counter, frame_feature in enumerate(x_lps):
8     ...

```

```

9     norm_feature = (frame_feature - mu) / sigma
10    frames.append(norm_feature)
11    x_lps = torch.stack(frames, dim=0)
12    ...

```

可以看到，NSNet的特征提取过程是将带噪语音经过STFT后，计算对数功率谱 $x\_lps$ ，将其进行标准化后作为网络的输入特征；幅度谱 $x\_ms$ 则用于与网络估计的时频掩蔽滤波得到增强的幅度谱。

网络部分则很简单，由三层GRU和一层全连接级联而成，最后的激活函数采用Sigmoid函数以保证网络的输出在[0,1]之间，这是由于IRM的范围即是如此(然而有很多工作表明将激活函数设为无界的ReLU可能带来性能的提升)。输入后的permute(.)操作是为了使变量维度变为(batch size, num frames, num bins)，使GRU和全连接对频率维进行操作。

```

1 # 网络部分 nsnet_model.py lines 37-53
2 def __build_model(self):
3     self.gru = nn.GRU(...)
4     self.dense = nn.Linear(...)
5
6     def forward(self, x):
7         x = x.permute(0, 2, 1)
8         x, _ = self.gru(x)
9         x = torch.sigmoid(self.dense(x))
10        x = x.permute(0, 2, 1)
11        return x

```

在训练过程中图1.4(也就是虚线下面的部分)，是作者提出的一种损失函数设计方式，同时计算了语音失真和噪声抑制：

$$\mathcal{L} = \alpha \mathbb{E}(\|S - \hat{M} \cdot S\|_2^2) + (1 - \alpha) \mathbb{E}(\|\hat{M} \cdot N\|_2^2), \quad (1.1)$$

可以看到，语音失真项损失是标准的信号近似(signal approximation, SA)形式，只用这一项作为损失函数的工作相对更多。要注意的是在训练阶段是使用纯净语音幅度谱与掩蔽的哈达玛积作为估计项的，有的工作则采用带噪幅度谱与掩蔽的积作为估计项，即 $\mathcal{L} = \mathbb{E}(\|S - \hat{M} \cdot X\|_2^2)$ 。式1.1对应的代码为：

```

1 # 损失部分 nsnet_model.py lines 57-71
2 def loss(self, target, prediction):
3     loss = F.mse_loss(prediction, target)
4     return loss
5
6 def training_step(self, batch, batch_idx):
7     ...
8     y_hat = self.forward(x_lps)
9     loss_speech = self.loss(y_ms [...], (y_hat * y_ms) [...])
10    loss_noise = self.loss(torch.zeros_like(y_hat), y_hat *
11                           noise_ms)
11    loss_val = self.alpha * loss_speech + (1 - self.alpha) *
12                           loss_noise
12    ...

```

此外，只计算掩蔽之间距离的掩蔽近似(mask approximation, MA)形式同样是基于掩蔽方法常见的损失函数，比如nngev的损失函数就是计算语音和噪声掩蔽之间的交叉熵损失(阅读nngev源码时请注意论文采用的是非因果网络BLSTM)：

```

1 # nn_models.py lines 17-22
2 def train_and_cv(self, Y, IBM_N, IBM_X, dropout=0.):
3     N_mask_hat, X_mask_hat = self._propagate(Y, dropout)
4     loss_X = binary_cross_entropy(X_mask_hat, IBM_X)
5     loss_N = binary_cross_entropy(N_mask_hat, IBM_N)
6     loss = (loss_X + loss_N) / 2
7     return loss

```

由于采用SA形式的损失函数，因此时频掩蔽并不需要显式定义。而采用时频掩蔽距离作为损失函数则需要显式计算出目标时频掩蔽，掩蔽的相关内容将在下一小节介绍。

接下来关注NSNet项目的解码过程，也就是测试阶段的代码。

```

1 # test_nn.py lines 27-34
2 gain_mask = model(x_lps)
3 y_spectrogram_hat = x_ms * gain_mask
4 y_stft_hat = torch.stack([y_spectrogram_hat * torch.cos(angle(
5     x_stft)), y_spectrogram_hat * torch.sin(angle(x_stft))], dim
6     =-1)
y_waveform_hat = istft(y_stft_hat, ...)

```

可以看到，基于掩蔽的语音增强算法在测试阶段模型将提取的带噪对数功率谱特征 $x_{lps}$ 作为输入并预测得到掩蔽 $gain\_mask$ ，而掩蔽与带噪幅度谱 $x_{ms}$ 进行哈达玛积得到滤波后的幅度谱 $y_{spectrogram\_hat}$ 。根据 $S_r = |S| \cdot \angle S$ 和 $S_i = |S| \cdot \angle S$ 的幅度相位-实部虚部关系将增强的幅度谱和带噪的相位耦合得到复数谱，并经过iSTFT重构回波形 $y_{waveform\_hat}$ 。

### 1.1.2 用于幅度谱的时频掩蔽

时频掩蔽的研究无疑是基于掩蔽的语音增强算法的重点。最早提出的时频掩蔽是理想二值掩蔽(Ideal Binary Mask, IBM)，根据语音时频分布的稀疏性，即语音成分在语谱中指分布于少数时频点，因此每个时频点上的语音和噪声之间的能量差异通常较大，于是可以将语音增强问题简单地看做将语音能量占主导的时频点筛选出来。二值掩蔽要做的，就是标记出那些语音占主导的时频点；而网络要做的，就是学习一个分类器，判断每个时频点是语音占主导还是噪声占主导。IBM被定义为：

$$M_{IBM}(t, f) = \begin{cases} 1, & \text{if } SNR(t, f) \geq T \\ 0, & \text{otherwise} \end{cases} \quad (1.2)$$

IBM的代码来自speech-feature-extractor:

```
1 # speech_utils.py lines 114-120
2 noise_spect = stft_extractor(noisy_speech-clean_speech, ...)
3 clean_spect = stft_extractor(clean_speech, ...)
4 ibm = np.where(10.*np.log10(clean_spect/noise_spect)>=local_snr
    , 1., 0.)
```

理想比值掩蔽(Ideal Ratio Mask, IRM)用一个0到1的实值代替了IBM的非0即1, 表征了时频点上语音能量和在语音噪声不相关假设下带噪语音能量的关系:

$$M_{IRM}(t, f) = \left( \frac{|S(t, f)|^2}{|S(t, f)|^2 + |N(t, f)|^2} \right)^{\beta} \quad (1.3)$$

式中 $\beta$ 一般取0.5。可以看出, IRM的定义式形式上与维纳滤波器完全一致, 但维纳滤波器是基于统计的线性滤波器, 式中的功率值是求期望的结果。代码同样来自speech-feature-extractor:

```
1 # speech_utils.py lines 122-129
2 noise_spect = stft_extractor(noisy_speech-clean_speech, ...)
3 clean_spect = stft_extractor(clean_speech, ...)
4 mask = np.sqrt(np.square(np.abs(clean_spect)) / (np.square(np.abs
    (noise_spect)) + np.square(np.abs(clean_spect))))
```

此外, 还有幅度谱掩蔽(Spectral Magnitude Mask, SMM)。相比IRM, SMM一方面在幅度上而非能量上进行定义, 另一方面也没有进行语音和噪声不相关的假设。这是由于根据IRM的定义, 在测试过程中的滤波过程 $\hat{M}_{IRM} \cdot |X| \neq |S|$ 。通过放弃语音和噪声不相关的假设, 将SMM定义为 $M_{SMM} = \frac{|S|}{|X|}$ , 其滤波后的结果是幅度谱上的最优滤波器估计, 然而, 这也导致了SMM的取值范围在 $[0, +\infty)$ 之间。早期为了使训练目标有界, SMM的最大值常被截断为1或2(代码中截断至2)。

```
1 # speech_utils.py line 10 and lines 130-136
2 EPSILON = np.finfo(np.float32).eps
3 ...
4 noisy_spect = stft_extractor(noisy_speech, ...)
5 clean_spect = stft_extractor(clean_speech, ...)
6 mask = np.abs(clean_spect) / (np.abs(noisy_spect) + EPSILON)
7 ...
8 mask = np.where(mask > 2., 2., mask)
```

相位敏感滤波器(Phase-sensitive mask, PSM)则是复数域上的最优的实值滤波器, 定义为:

$$M_{PSM} = \frac{|S|}{|X|} \cos(\angle S - \angle X), \quad (1.4)$$

可以看出, PSM的范围是 $(-\infty, +\infty)$ 。和SMM一样, PSM也常被截断以方便训练, 常用的截断范围为 $[0, 1]$ 。

若采用SA的方式训练，损失函数定义为：

$$\mathcal{L} = \mathbb{E}\{(\hat{M}_{PSM} \cdot |X| - |S| \cos(\angle S - \angle X))^2\}, \quad (1.5)$$

而解码时和上述掩蔽一样直接作用于带噪幅度谱后与带噪相位耦合得到增强的谱。

PSM的代码选自espnet，代码中的 $r$ 是纯净语音的复谱，利用 $\angle S = \frac{s}{|S|}$ 的方式得到相位，并利用三角函数关系式得到相位差的余弦值：

```

1 # espnet2/enh/loss/criterions/tf_domain.py lines 63-73
2 phase_r = r / (abs(r) + EPS)
3 phase_mix = mix_spec / (abs(mix_spec) + EPS)
4 # cos(a - b) = cos(a)*cos(b) + sin(a)*sin(b)
5 cos_theta = phase_r.real * phase_mix.real + phase_r.imag *
   phase_mix.imag
6 mask = (abs(r) / (abs(mix_spec) + EPS)) * cos_theta
7 mask = (
8     mask.clamp(min=0, max=1)
9     ...
10 )

```

## 1.2 基于谱映射的语音增强算法

### 1.2.1 处理流程

基于谱映射的语音增强算法相比之下对前置知识的要求更低，与基于深度学习的图像分割任务类似，这类方法将语谱视为一张图片，网络接收带噪语谱图作为输入，直接预测得到增强的语谱图。本节将以CRN-causal为例展示基于谱映射的语音增强算法的处理过程。

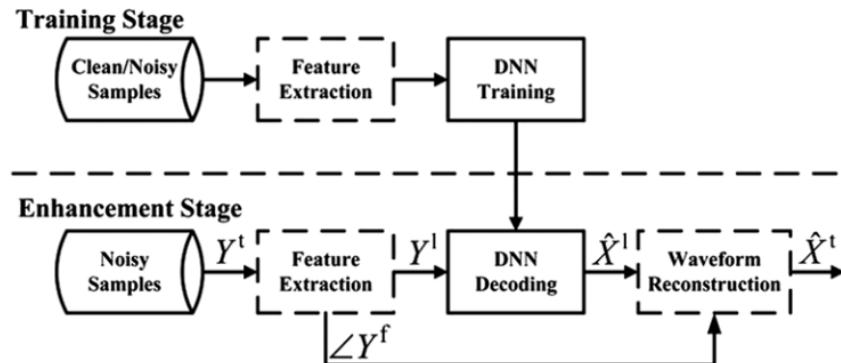


Fig. 1.5 基于谱映射的幅度谱算法框图[5]

尽管图1.5给出的是SEDNN的框图，但是我们将结合CRN-causal说明谱映射方法的流程。这张流程图看起来就简明得多，和图1.4相反，虚线以上是训练过程而虚线以下是解码过程。训练阶段的带噪语音和纯净语音首先经过特征提取，而后将特征送入DNN。经过损失函数训练后得到优化后的网络参数。测试阶段则将带噪语音 $Y^t$ 经过特征提取得到幅度谱 $Y^1$ (而SEDNN的特征是对数功率谱)和相位谱 $ZY^f$ 。幅度谱 $Y^1$ 送入训练好的DNN得到增强的幅度谱 $\hat{X}^1$ 。最后经过波形重构得到增强的时域波形 $\hat{X}^t$ 。

CRN-causal直接以STFT幅度谱作为输入特征和映射目标。为了方便转成C语言落地，这里的STFT采用以傅里叶基为卷积核的一维卷积的方式，而后利用 $|X| = \sqrt{\mathcal{R}(X)^2 + \mathcal{I}(X)^2}$ 计算得到幅度谱，具体代码如下：

```

1 # scripts/utils/pipeline_modules.py lines 7-21
2 class NetFeeder(object):
3     def __init__(self, device, win_size=320, hop_size=160):
4         self.eps = torch.finfo(torch.float32).eps
5         self.stft = STFT(win_size, hop_size).to(device)
6
7     def __call__(self, mix, sph):
8         real_mix, imag_mix = self.stft.stft(mix)
9         mag_mix = torch.sqrt(real_mix**2 + imag_mix**2)
10        feat = mag_mix
11
12        real_sph, imag_sph = self.stft.stft(sph)
13        mag_sph = torch.sqrt(real_sph**2 + imag_sph**2)
14        lbl = mag_sph
15        return feat, lbl

```

代码中 $feat$ 和 $lbl$ 分别作为网络的特征和训练目标用于对网络进行训练：

```

1 # scripts/utils/models.py lines 162-174
2 # forward + backward + optimize
3 optimizer.zero_grad()
4 with torch.enable_grad():
5     est = net(feat)
6     loss = criterion(est, lbl, loss_mask, n_frames)
7     loss.backward()
8     if self.clip_norm >= 0.0:
9         clip_grad_norm_(net.parameters(), self.clip_norm)
10        optimizer.step()
11    # calculate loss
12    running_loss = loss.data.item()
13    accu_tr_loss += running_loss * sum(n_frames)
14    accu_n_frames += sum(n_frames)

```

$feat$ 输入网络得到估计的谱 $est$ ，而后计算 $est$ 和 $lbl$ 之间的损失反向传播优化网络参数。

解码过程同样采用NetFeeder类的实例feeder进行特征提取，将幅度谱特征 $feat$ 送入训练好的网络得到估计的幅度谱 $est$ ，最后经过波形合成过程resynthesizer即可得到增强语音：

```

1 # scripts/utils/models.py lines 337-348
2 feat, lbl = feeder(mix, sph)

```

```

3   with torch.no_grad():
4     ...
5     est = net(feats)
6     ...
7   ...
8   sph_est = resynthesizer(est, mix)

```

波形合成resynthesizer对应的类Resynthesizer按照

$$\begin{aligned}\mathcal{R}(\hat{S}) &= |\hat{S}| \cdot \cos(\angle X), \\ \mathcal{I}(\hat{S}) &= |\hat{S}| \cdot \sin(\angle X)\end{aligned}\quad (1.6)$$

得到增强的复谱，经过iSTFT生成增强的语音波形：

```

1 # scripts/utils/models.py lines 337-348
2 class Resynthesizer(object):
3   def __init__(self, device, win_size=320, hop_size=160):
4     self.stft = STFT(win_size, hop_size).to(device)
5
6   def __call__(self, est, mix):
7     real_mix, imag_mix = self.stft.stft(mix)
8     pha_mix = torch.atan2(imag_mix.data, real_mix.data)
9     real_est = est * torch.cos(pha_mix)
10    imag_est = est * torch.sin(pha_mix)
11    sph_est = self.stft.istft(torch.stack([real_est, imag_est],
12                                          dim=1))
13    sph_est = F.pad(sph_est, [0, mix.shape[1]-sph_est.shape[1]])
14
15   return sph_est

```

最后简单提一下网络结构，网络结构如图1.6所示，结合图1.6和代码可以看出，网络可以分为三部分：编码器、时序建模结构和解码器。编码器由5层二维卷积级联Batch Normalization和ReLU激活函数组成；时序建模结构是两层LSTM；解码器由5层二维反卷积构成，除了最后一层反卷积层，其余反卷积层同样级联了Batch Normalization和ReLU激活函数，最后一层反卷积层后则是Softplus激活函数以保证输出结果为非负数从而满足幅度谱的物理意义。幅度谱经过编码器被变换到高维嵌入式(embedding)空间，卷积核用于建模相邻帧和相邻频点之间的相关性。而后通道维和频率特征维被合并到一个维度，经过两层LSTM建模不同帧之间这组特征的时序关系，最后变回与编码器输出相同的张量形状，经过解码器张成与输入幅度谱相同的尺寸，从而得到增强的幅度谱。编解码器之间采用concatenate形式的skip connection连接以缓解网络过深可能导致的梯度消失问题。可以看出，该网络架构与图像分割中的编解码器结构极其相似，这是不难理解的，因为以上两个任务都可以看做逐点(时频点/像素点)的滤波问题。然而二者仍有不同之处，相比于图片是一次快拍(借用阵列信号处理中的概念，可以认为是同一时刻)即可全部获得的，语谱图由于其中一个维度是时间帧，需要等待一句话结束后才能得到一张语谱图。在实时应用的推理阶段每次只能获得当前帧而未来帧的信息是未知的，因此需要一些技术来保证模型的因果性。这里的代码在编解码器中

沿时间帧维度的截断和补零被使用以保证卷积操作的因果性，而我们将在第3章中对因果性进行更详细的讨论。

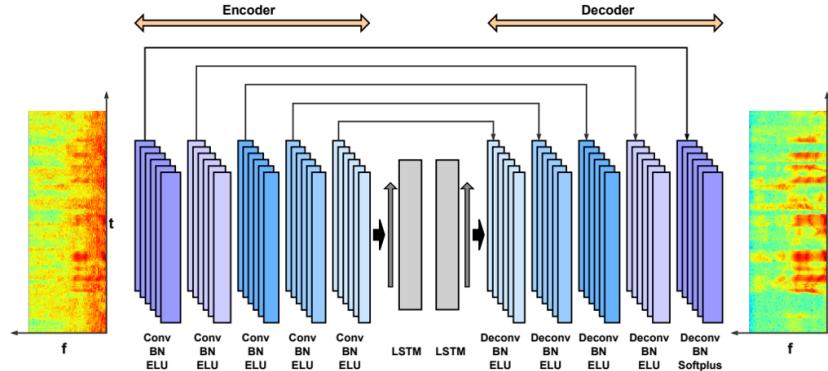


Fig. 1.6 CRN-causal 框图[6]

```

1 # scripts/utils/networks.py
2 class Net(nn.Module):
3     def __init__(self):
4         super(Net, self).__init__()
5
6         self.conv1 = nn.Conv2d(in_channels=1, out_channels=16,
7             kernel_size=(2,3), stride=(1,2), padding=(1,0))
8         self.conv2 = nn.Conv2d(in_channels=16, out_channels=32,
9             kernel_size=(2,3), stride=(1,2), padding=(1,0))
10        self.conv3 = nn.Conv2d(in_channels=32, out_channels=64,
11            kernel_size=(2,3), stride=(1,2), padding=(1,0))
12        self.conv4 = nn.Conv2d(in_channels=64, out_channels=128,
13            kernel_size=(2,3), stride=(1,2), padding=(1,0))
14        self.conv5 = nn.Conv2d(in_channels=128, out_channels=256,
15            kernel_size=(2,3), stride=(1,2), padding=(1,0))
16
17        self.lstm = nn.LSTM(256*4, 256*4, 2, batch_first=True)
18
19        self.conv5_t = nn.ConvTranspose2d(in_channels=512,
20            out_channels=128, kernel_size=(2,3), stride=(1,2),
21            padding=(1,0))
22        self.conv4_t = nn.ConvTranspose2d(in_channels=256,
23            out_channels=64, kernel_size=(2,3), stride=(1,2),
24            padding=(1,0))
25        self.conv3_t = nn.ConvTranspose2d(in_channels=128,
26            out_channels=32, kernel_size=(2,3), stride=(1,2),
27            padding=(1,0))
28        self.conv2_t = nn.ConvTranspose2d(in_channels=64,
29            out_channels=16, kernel_size=(2,3), stride=(1,2),
30            padding=(1,0), output_padding=(0,1))

```

```

18     self.conv1_t = nn.ConvTranspose2d(in_channels=32,
19         out_channels=1, kernel_size=(2,3), stride=(1,2),
20         padding=(1,0))
21
22     self.bn1 = nn.BatchNorm2d(16)
23     self.bn2 = nn.BatchNorm2d(32)
24     self.bn3 = nn.BatchNorm2d(64)
25     self.bn4 = nn.BatchNorm2d(128)
26     self.bn5 = nn.BatchNorm2d(256)
27
28     self.bn5_t = nn.BatchNorm2d(128)
29     self.bn4_t = nn.BatchNorm2d(64)
30     self.bn3_t = nn.BatchNorm2d(32)
31     self.bn2_t = nn.BatchNorm2d(16)
32     self.bn1_t = nn.BatchNorm2d(1)
33
34     self.elu = nn.ELU(inplace=True)
35     self.softplus = nn.Softplus()
36
37     def forward(self, x):
38
39         out = x.unsqueeze(dim=1)
40         e1 = self.elu(self.bn1(self.conv1(out)[:, :, :-1, :].
41             contiguous()))
42         e2 = self.elu(self.bn2(self.conv2(e1)[:, :, :-1, :].
43             contiguous()))
44         e3 = self.elu(self.bn3(self.conv3(e2)[:, :, :-1, :].
45             contiguous()))
46         e4 = self.elu(self.bn4(self.conv4(e3)[:, :, :-1, :].
47             contiguous()))
48         e5 = self.elu(self.bn5(self.conv5(e4)[:, :, :-1, :].
49             contiguous()))
50
51         out = e5.contiguous().transpose(1, 2)
52         q1 = out.size(2)
53         q2 = out.size(3)
54         out = out.contiguous().view(out.size(0), out.size(1), -1)
55         out, _ = self.lstm(out)
56         out = out.contiguous().view(out.size(0), out.size(1), q1, q2)
57         out = out.contiguous().transpose(1, 2)
58
59         out = torch.cat([out, e5], dim=1)
60
61         d5 = self.elu(torch.cat([self.bn5_t(F.pad(self.conv5_t(
62             out), [0, 0, 1, 0]).contiguous(), e4], dim=1)))
63         d4 = self.elu(torch.cat([self.bn4_t(F.pad(self.conv4_t(d5
64             ), [0, 0, 1, 0]).contiguous(), e3], dim=1)))
65         d3 = self.elu(torch.cat([self.bn3_t(F.pad(self.conv3_t(d4
66             ), [0, 0, 1, 0]).contiguous(), e2], dim=1)))
67         d2 = self.elu(torch.cat([self.bn2_t(F.pad(self.conv2_t(d3
68             ), [0, 0, 1, 0]).contiguous(), e1], dim=1)))
69         d1 = self.softplus(self.bn1_t(F.pad(self.conv1_t(d2),
70             [0, 0, 1, 0]).contiguous())))

```

```

60     out = torch.squeeze(d1, dim=1)
61
62     return out

```

### 1.3 小结

本章介绍了基于幅度谱深度学习语音增强技术的两种主流方法：基于掩蔽的语音增强和基于谱映射的语音增强。这两类算法是后续语音增强模型的基础，并将在后续章节中进一步发展。

## References

1. 刘文举等: 基于深度学习语音分离技术的研究现状与进展, 自动化学报, 42(6), pp. 819-833, 2016.
2. Wang, et, al. *supervised speech separation based on deep learning: an overview*, IEEE/ACM Transactions on Audio, Speech, and Language Processing, 26(10), pp.1702-1726, 2018.
3. Y. Xia, et. al. *Weighted Speech Distortion Losses for Neural-network-based Real-time Speech Enhancement*, ICASSP 2020.
4. X. Hao, et. al. *Masking and Inpainting: A Two-Stage Speech Enhancement Approach for Low SNR and Non-Stationary Noise*, ICASSP 2020, pp. 6959–6963. doi: 10.1109/I-CASSP40776.2020.9053188.
5. Y. Xu, et. al. *An Experimental Study on Speech Enhancement BasedonDeepNeuralNetworks*, IEEE Signal Processing Letters, 21(1), pp. 65-68, 2014.
6. K. Tan, et. al. *A convolutional recurrent neural network for real-time speech enhancement*, INTERSPEECH 2018, pp. 3229-3233.
7. S. Braun, et. al. *Effect of noise suppression losses on speech distortion and ASR performance*, arXiv:2111.11606.
8. N. Zheng, et. al. *Phase-Aware Speech Enhancement Based on Deep Neural Networks*, IEEE/ACM Transactions on Audio, Speech, and Language Processing, 27(1), pp.63-76, 2019.
9. K. Tan, et. al. *Learning complex spectral mapping with gated convolutional recurrent networks for monaural speech enhancement*, 28, pp.380-390, 2020.
10. H. Choi, et. al. *Phase-aware speech enhancement with deep complex U-Net*, arXiv:1903.03107.
11. X. Le, et. al. *DPCRN: Dual-Path Convolution Recurrent Network for Single Channel Speech Enhancement*, Interspeech 2021, pp. 2811-2815.
12. S. Braun, et. al. *Towards efficient models for real-time deep noise suppression*, arXiv:2101.09249.
13. D. Yin, et. al. *PHASEN: a phase-and-harmonics-aware speech enhancement network*, Proceedings of the AAAI Conference on Artificial Intelligence, 34(5), pp. 9458-9465, 2020.
14. Z. Wang, et. al. *Deep Learning Based Phase Reconstruction for Speaker Separation: A Trigonometric Perspective*, ICASSP 2019, pp. 71-75.
15. S. Suhadi, et. al. *A data-driven approach to a priori SNR estimation*, IEEE Transactions on Audio, Speech and Language Processing, vol. 19, no. 1, pp. 186–195, 2011.
16. Y. Ephraim and D. Malah. *Speech enhancement using a minimum mean-square error log-spectral amplitude estimator*, IEEE Trans. Acoust. Speech Signal Process. 33 (2), pp. 443–445. 1985.
17. Y. Xia, et. al. *Low-dimensional recurrent neural network-based Kalman filter for speech enhancement*, Neural Networks, 67, pp. 131-139, 2015.

## Chapter 2

### 基于复数谱的语音增强

随着深度学习的应用和发展，基于幅度谱的语音增强算法相比传统算法取得了显著的性能提升。然而，语谱包含幅度和相位两个分量。这类算法是对幅度进行修正而仍使用带噪信号的相位谱显然不够完美。面对当时语音增强算法的性能瓶颈，人们难免不会将此归咎于在相位谱上的不作为。相位谱对语音质量恢复的重要性无疑为当时的人们提供了一个绝佳的佐证。因此，一系列意图恢复相位的工作不断被提出。

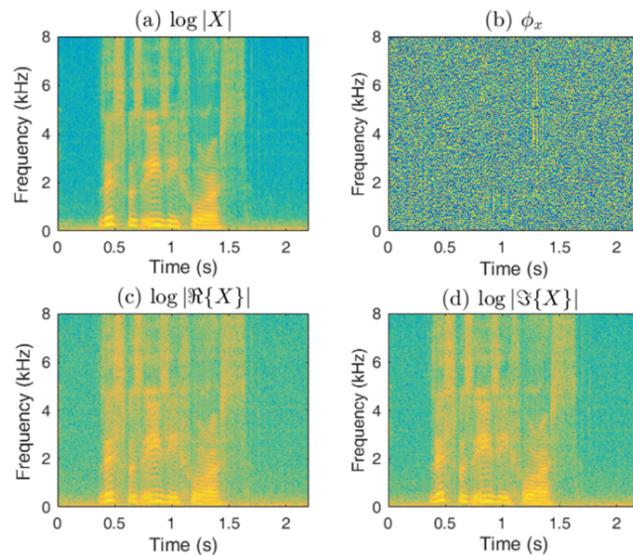


Fig. 2.1 语音复谱图[8]

在谈及基于复数谱语音增强方法的分类之前，有必要对相位谱估计的困难性进行说明。图2.1(b)展示了一段语音的相位谱。可以看出，由于各时频点的相位在时间轴和频率轴上都变化迅速，由于相位缠绕问题，语谱各时频点

的相位几乎均匀分布在 $[-\pi, \pi)$ 之间。相位谱并不存在一个清晰的结构模式，导致网络很难对其学习预测。而图2.1(c)和(d)展示的复谱的实部和虚部则具有与幅度谱类似的结构模式，使得通过网络估计实部虚部分量从而隐式优化幅度和相位成为可能。

于是，目前针基于复数谱的语音增强方法大致可分为三类：构造值数掩蔽、直接复数谱映射和估计幅度相位分量，其中复值掩蔽按坐标系的区别可分为笛卡尔坐标系下的时频掩蔽和极坐标系下的时频掩蔽。复值掩蔽和复数谱映射通过神经网络预测掩蔽或复谱的实部和虚部。而幅度相位分量估计则是在原有的(基于掩蔽或谱映射的)幅度谱语音增强技术的基础上额外构造模块估计相位谱或与相位相关的群时延、瞬时频率导数等分量。

尽管我们可以大致确认重构相位谱的任务对目前的深度学习技术是影响模型鲁棒性甚至是优化困难的，但将其归因于相位谱缺乏清晰的结构模式似乎只是一种假设解释。另外，需要注意的是试验只表明了相位谱的估计是困难的，但相位谱在一些语音处理领域作为输入特征是成功的。

本章安排如下，在2.1中是基于复数谱映射的语音增强算法的相关描述，基于笛卡尔坐标系下的复值掩蔽和基于极坐标系下的复值掩蔽的介绍分别在2.2和2.3，最后的2.4介绍了幅度相位估计算法的相关工作。

## 2.1 基于复数谱映射的语音增强

基于复数谱的语音增强模型，顾名思义，就是将复数谱的实部和虚部成分作为网络输入和输出。实现起来也可以非常简单，只需将基于幅度谱的语音增强模型的输入从一个通道的幅度谱改为两个通道的复数谱(实部和虚部)，并将网络的输出从利用激活函数保证非负性的单通道的幅度变成不再使用任何激活函数的双通道的复数谱即可。一个合理的疑问可能会出现在大多数读者心中：直接将实部和虚部沿特征通道拼接是否是合理的，换言之，实部和虚部成分是否应该单独的编码或预测。Tan在2019年通过实验比较了这几种网络架构的性能区别，不出意外地以上结构都可以实现复数谱映射，而Tan更推荐使用单独的编码器同时编码带噪谱的实部和虚部，而后使用两个解码器分别估计增强语谱的实部和虚部。该模型后来被成为GCRN，其卓越的性能成为了后续很多语音增强工作的基础。因此，这里也采用GCRN作为基于复数谱映射语音增强模型的例子，通过其代码理解这类算法。

GCRN的模型框架和网络参数如图2.2所示。该工作延续了Tan之前的CRN-causal(见1.2)的因果卷积编解码器结合LSTM时序建模的架构。带噪语谱的实部和虚部沿特征通道维拼接后送入一个编码器，编码后的高维特征经过LSTM进行建模后作为两个解码器的输入，两个解码器分别估计得到增强语谱的实部和虚部。其中，卷积和反卷积层用其门控(gating)结构代替以提升语音增强性能，LSTM也引入了分组(grouped)机制减少模型的参数量和计算

复杂度，这两个改进点至今仍被许多工作所采用。GCRN的模型前向处理流程为：

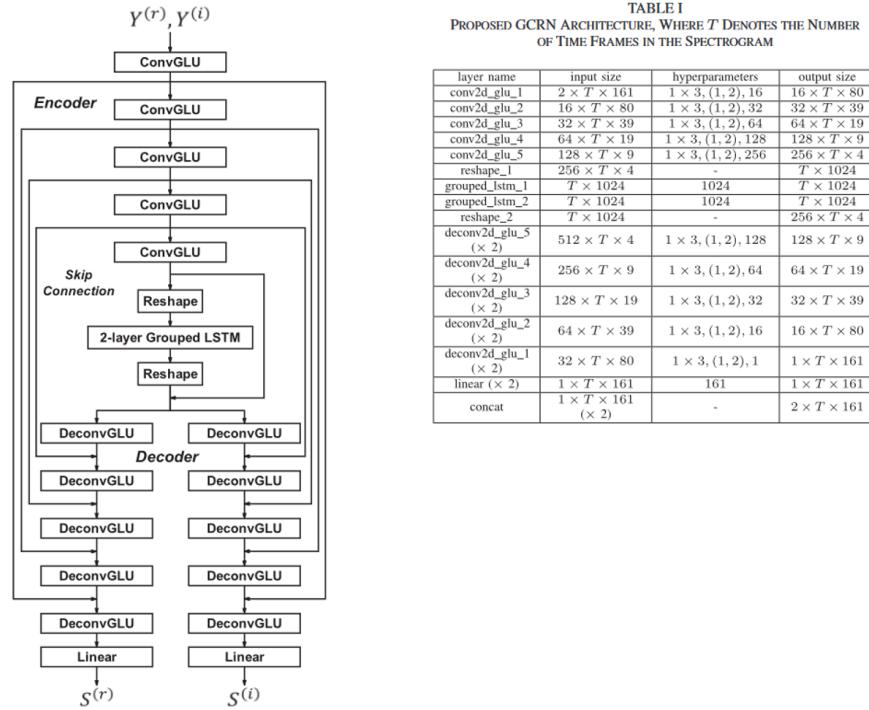


Fig. 2.2 GCRN框架[9]

```

1      # utils/networks.py lines 134-165
2      def forward(self, x):
3          out = x
4          e1 = self.elu(self.bn1(self.conv1(out)))
5          e2 = self.elu(self.bn2(self.conv2(e1)))
6          e3 = self.elu(self.bn3(self.conv3(e2)))
7          e4 = self.elu(self.bn4(self.conv4(e3)))
8          e5 = self.elu(self.bn5(self.conv5(e4)))
9          out = e5
10         out = self.glstm(out)
11         out = torch.cat((out, e5), dim=1)
12
13         d5_1 = self.elu(torch.cat((self.bn5_t_1(self.conv5_t_1(
14             out)), e4), dim=1))
15         d4_1 = self.elu(torch.cat((self.bn4_t_1(self.conv4_t_1(
16             d5_1)), e3), dim=1))
17         d3_1 = self.elu(torch.cat((self.bn3_t_1(self.conv3_t_1(
18             d4_1)), e2), dim=1))

```

```

16     d2_1 = self.elu(torch.cat((self.bn2_t_1(self.conv2_t_1(
17         d3_1)), e1), dim=1))
18     d1_1 = self.elu(self.bn1_t_1(self.conv1_t_1(d2_1)))
19
20     d5_2 = self.elu(torch.cat((self.bn5_t_2(self.conv5_t_2(
21         out)), e4), dim=1))
22     d4_2 = self.elu(torch.cat((self.bn4_t_2(self.conv4_t_2(
23         d5_2)), e3), dim=1))
24     d3_2 = self.elu(torch.cat((self.bn3_t_2(self.conv3_t_2(
25         d4_2)), e2), dim=1))
26     d2_2 = self.elu(torch.cat((self.bn2_t_2(self.conv2_t_2(
27         d3_2)), e1), dim=1))
28     d1_2 = self.elu(self.bn1_t_2(self.conv1_t_2(d2_2)))
29
30     return out

```

GCRN的训练流程如下:

```

1 # utils/models.py lines 145-174
2 for n_iter, egs in enumerate(tr_loader):
3     n_iter += start_iter
4     mix = egs['mix']
5     sph = egs['sph']
6     n_samples = egs['n_samples']
7
8     mix = mix.to(self.device)
9     sph = sph.to(self.device)
10    n_samples = n_samples.to(self.device)
11    n_frames = countFrames(n_samples, self.win_size, self.
12                           hop_size)
13    ...
14    # prepare features and labels
15    feat, lbl = feeder(mix, sph)
16    loss_mask = lossMask(shape=lbl.shape, n_frames=n_frames,
17                          device=self.device)
18    # forward + backward + optimize
19    optimizer.zero_grad()
20    with torch.enable_grad():
21        est = net(feat)
22        loss = criterion(est, lbl, loss_mask, n_frames)
23        loss.backward()
24        if self.clip_norm >= 0.0:
25            clip_grad_norm_(net.parameters(), self.clip_norm)
26        optimizer.step()
27    # calculate loss
28    running_loss = loss.data.item()
29    accu_tr_loss += running_loss * sum(n_frames)
30    accu_n_frames += sum(n_frames)

```

和CRN-causal一样，feeder是用于特征提取的对象，提取了带噪语音和纯净语音的复数谱。

## 2.2 基于笛卡尔坐标系的复值掩蔽方法

和基于幅度谱的语音增强方法类似，基于复数谱的语音增强方法除了基于谱映射的形式之外也包括使用掩蔽滤波的方式。网络的输出结果不再是复数谱的实部和虚部成分而是时频掩蔽。复值理想比值掩蔽(complex-valued ideal ratio mask, cIRM)于2016年被提出，处理过程如图2.3-1所示。它将复数域分解为实部和虚部的组合形式，即 $M = M_r + j \cdot M_i \in \mathbb{C}$ 。利用神经网络直接估计实部掩蔽和虚部掩蔽，而后将估计的掩蔽与带噪复数谱按复数乘法相乘进行滤波得到增强语谱：

$$\begin{aligned}\widehat{S}_r &= \widehat{M}_r \cdot X_r - \widehat{M}_i \cdot X_i, \\ \widehat{S}_i &= \widehat{M}_r \cdot X_i + \widehat{M}_i \cdot X_r,\end{aligned}\quad (2.1)$$

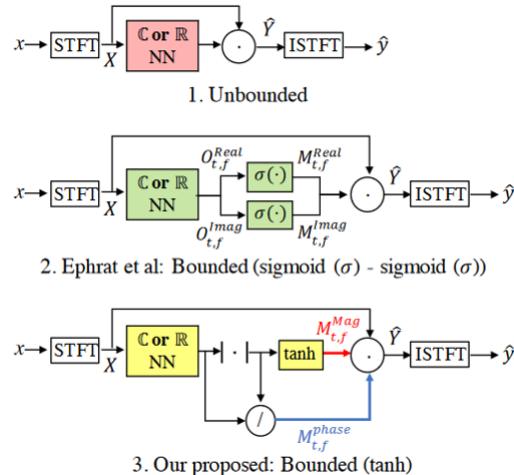


Fig. 2.3 基于复数谱的掩蔽示意图[10]

如果采用MA的方式计算损失函数，则需要按如下方式计算cIRM作为标签：

$$\begin{aligned}M_r &= \frac{X_r \cdot S_r + X_i \cdot S_i}{X_r^2 + X_i^2}, \\ M_i &= \frac{X_r \cdot S_i - X_i \cdot S_r}{X_r^2 + X_i^2}\end{aligned}\quad (2.2)$$

为了与2.3中介绍的掩蔽区分，有时也将cIRM成为基于笛卡尔坐标系的复值掩蔽。需要注意的是，为了便于网络优化，cIRM最初采用有界的形式通过tanh函数将范围压缩至某个对称的范围内。然后后续的工作表明至少采用SA的方式采用无界的cIRM并没有出现优化问题，DPCRN展示了掩蔽的滤波过程：

```

1 # main.py lines 232-246
2 def mk_mask(self, x):
3     ...
4     enh_real = noisy_real * mask_real - noisy_imag *
5         mask_imag
6     enh_imag = noisy_real * mask_imag + noisy_imag *
7         mask_real
8
9     return [enh_real, enh_imag]

```

本节余下部分将简要介绍DPCRN的源码，首先DPCRN的框架如图2.4所示 可以看到，与GCRN相似，网络架构由编码器、DPRNN构成的时序建

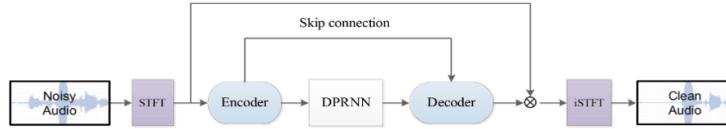


Fig. 2.4 DPCRN框架[11]

模模块和解码器组成，编解码器与GCRN相似——五层由二维因果卷积层+BN层+PReLU层构成的编码层和五层由二维因果转置卷积+BN层+PReLU层构成的解码层(为了保证输出无界，最后一个转置卷积不级联BN层和PReLU)，只是舍弃了门控卷积并且只使用了一个解码器。时序建模模块的DPRNN源自[]，为了发挥DPRNN的作用，编码器输出的特征维度设为较大的128。编码器的输入维度设为2以使带噪复数谱按实部、虚部沿特征通道维度拼接后输入，对应的，解码器的输出维度也设为2从而得到cIRM的实部和虚部。值得注意的是一种instant layer normalization(iLN)以应对输入语音动态范围过大问题。DPCRN在谱域上采用SA的损失函数未对掩蔽约束至有界。模型代码如下：

```

1 # main.py lines 248-341
2 def build_DPCRN_model(self, ...):
3     ...
4     '''encoder'''
5     ...
6     input_complex_spec = LayerNormalization(...)(
7         input_complex_spec)
8     ...
9     conv_1 = Conv2D(32, (2,5),(1,2),...)(input_complex_spec)
10    bn_1 = BatchNormalization(...)(conv_1)
11    out_1 = PReLU(shared_axes=[1,2])(bn_1)
12    ...

```

```

12 conv_5 = Conv2D(128, (2,3),(1,1),...)(out_4)
13 bn_5 = BatchNormalization(...)(conv_5)
14 out_5 = PReLU(shared_axes=[1,2])(bn_5)
15
16 dp_in = out_5
17 for i in range(self.numDP):
18     dp_in = DprnnBlock(...)(dp_in)
19     dp_out = dp_in
20
21     '''decoder'''
22 skipcon_1 = Concatenate(axis = -1)([out_5,dp_out])
23 deconv_1 = Conv2DTranspose(64,(2,3),(1,1),...)(skipcon_1)
24 dbn_1 = BatchNormalization(...)(deconv_1)
25 dout_1 = PReLU(shared_axes=[1,2])(dbn_1)
26 ...
27 skipcon_5 = Concatenate(axis = -1)([out_1,dout_4])
28 deconv_5 = Conv2DTranspose(2,(2,5),(1,2),...)(skipcon_5)
29 ...
30 output_mask = deconv_5
31
32 enh_spec = Lambda(self.mk_mask)([real,imag,output_mask])
33 ...

```

除此之外，这里不得不提及S. Braun等人于2021年提出的CRUSE模型，其作为一个参数量和计算量较低的实时语音增强模型在ICASSP2021的DNS Challenge中取得了极其优异的成绩。网络架构如图2.5所示，同样是编解码器结构，和DPCRN一样(但它比DPCRN更早提出)，门控卷积是被舍弃的，解码器也减至一个。此外GCRN的分组机制被保留只是RNN从LSTM变为计算量和参数量更少的GRU，skip connection也不再以拼接的方式而是以相加的方式(如图2.6)将编码器的信息送入解码器中，进一步降低了参数量和计算量。

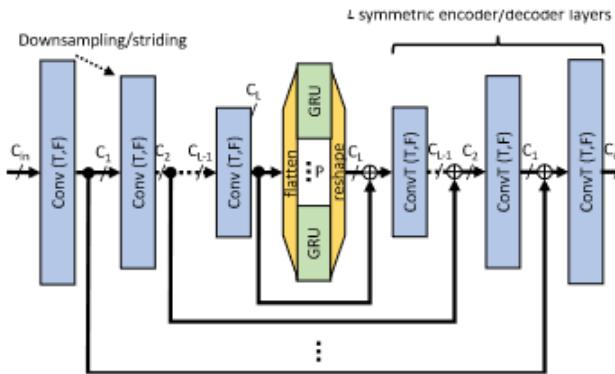


Fig. 2.5 CRUSE框架[12]

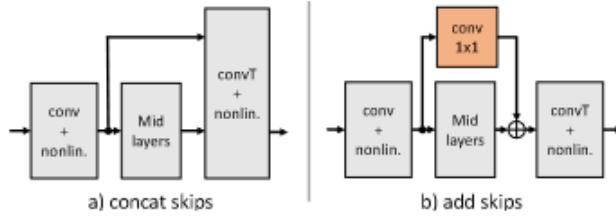


Fig. 2.6 skip connection的改进[12]

CRUSE的特征采用了幂律压缩的复数谱( $X_{cprs} = \frac{X(k,n)}{|X(k,n)|} |X(k,n)|^{0.3}$ )，网络输出STFT域的cIRM，对带噪复谱进行滤波后采用能量级不变的幂律压缩RI+Mag损失函数结合STFT一致性约束进行网络优化。以上训练流程如图2.7所示损失函数定义为

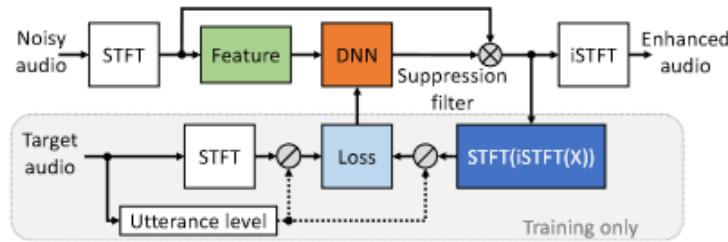


Fig. 2.7 CRUSE训练流程[12]

$$\mathcal{L} = \frac{1}{\sigma_S^c} (\lambda \sum_{k,n} |S^c - \hat{S}^c|^2 + (1 - \lambda) (\sum_{k,n} ||S|^c - |\hat{S}|^c|^2)), \quad (2.3)$$

其中 $\sigma_S$ 是纯净语音有声段的能量， $c$ 设为0.3。

处理针对当前时频点的复值掩蔽，利用当前时频点相邻时频点进行滤波求和计算得到当前时频点谱估计的深度滤波(deep filtering)方法近年来也开始被使用。深度滤波的形式可写作：

$$\hat{S}(t, f) = \sum_{k=-K}^K \sum_{i=0}^N \hat{G}(l, k) \odot X(t - i + l, f - k). \quad (2.4)$$

式中 $K$ 和 $N$ 分别为与当前时频点相邻的频点数和帧数， $l$ 为前看(look-ahead)偏移，若为因果模型则 $l$ 取0。有时 $K$ 取0，即指利用相同频率的当前帧和历史帧信息进行滤波，而不使用相邻频点信息。 $\odot$ 是指复数乘法。

## 2.3 基于极坐标系的复值掩蔽方法

区别于基于笛卡尔坐标系下的复值掩蔽，基于极坐标系的复值掩蔽将掩蔽对带噪复谱的滤波作用按如下方式理解：

$$\widehat{S} = \widehat{M} \cdot X = |\widehat{M}| \cdot X \cdot e^{j(\angle \widehat{M} + \angle X)} \quad (2.5)$$

复值掩蔽可视为对每个时频点幅度的放缩和相位的旋转。为了方便网络优化，作用在幅度谱上的掩蔽通过  $\tanh(\cdot)$  函数将无界范围非线性映射到复平面的单位圆内。基于极坐标系的复值掩蔽计算方法如图2.3(b)-3所示，网络隐式的估计该掩蔽的实部和虚部成分后，分别计算得到无界的幅度掩蔽分量和相位掩蔽分量，无界的幅度掩蔽分量通过  $\tanh(\cdot)$  激活函数得到有界的幅度掩蔽分量。而后将估计得到的掩蔽按式2.5对带噪复谱进行滤波，从而得到增强语音的复谱。H. S. Choi等通过实验分析了基于笛卡尔坐标系和极坐标系的复值掩蔽在谱域和时域损失函数下处理的结果，通过波形对比表明了基于极坐标系的复值掩蔽在时域损失函数下处理结果与纯净语音波形更相似，而谱域损失函数并不利于相位重构。不过需要注意的是，该结论是基于波形相似度得出的，在听感上该结论是否成立是还需验证的。

这里采用DCCRN的代码进行举例。该模型输入是将实部和虚部沿特征通道维度拼接的带噪的复数谱，网络输出的是基于极坐标系的复值掩蔽，而后复值掩蔽对带噪语音的幅度和相位谱进行调整，并通过iSTFT变换回时域的增强波形。其中网络采用了编码器-时序建模-解码器架构，与之前提出方法不同是，各模块都采用了复值网络。复值网络简单来说就是将特征沿特征通道维对半分成“实部” $E_r$ 和“虚部” $E_i$ ，二者分别经过两个参数不共享的网络层  $f_r(\cdot)$  和  $f_i(\cdot)$ ，得到输出的  $f_r(E_r)$ 、 $f_i(E_r)$ 、 $f_r(E_i)$  和  $f_i(E_i)$ ，而后这四个分量模仿复数乘法的计算规则得到最终结果  $(f_r(E_r) - f_i(E_i)) + j(f_r(E_i) + f_i(E_r))$ 。

```

1 # dc_crn.py lines 150-234
2 def forward(self, inputs, lens=None):
3     # 特征提取部分
4     specs = self.stft(inputs)
5     real = specs[:, :self.fft_len // 2 + 1]
6     imag = specs[:, self.fft_len // 2 + 1:]
7     spec_mags = torch.sqrt(real ** 2 + imag ** 2 + 1e-8)
8     spec_phase = torch.atan2(imag, real)
9     cspecs = torch.stack([real, imag], 1)
10    cspecs = cspecs[:, :, 1:]
11    out = cspecs
12
13    # 网络部分
14    encoder_out = []
15    for idx, layer in enumerate(self.encoder):
16        out = layer(out)
17        encoder_out.append(out)
18
19    batch_size, channels, dims, lengths = out.size()
20    out = out.permute(3, 0, 1, 2)
21    if self.use_clstm:

```

```

22     r_rnn_in = out[:, :, : channels // 2]
23     i_rnn_in = out[:, :, channels // 2:]
24     r_rnn_in = torch.reshape(r_rnn_in, [lengths, batch_size,
25                               channels // 2 * dims])
26     i_rnn_in = torch.reshape(i_rnn_in, [lengths, batch_size,
27                               channels // 2 * dims])
28     r_rnn_in, i_rnn_in = self.enhance([r_rnn_in, i_rnn_in])
29     r_rnn_in = torch.reshape(r_rnn_in, [lengths, batch_size,
30                               channels // 2, dims])
31     i_rnn_in = torch.reshape(i_rnn_in, [lengths, batch_size,
32                               channels // 2, dims])
33     out = torch.cat([r_rnn_in, i_rnn_in], 2)
34     ...
35     out = out.permute(1, 2, 3, 0)
36
37     for idx in range(len(self.decoder)):
38         out = complex_cat([out, encoder_out[-1 - idx]], 1)
39         out = self.decoder[idx](out)
40         out = out[... , 1:]
41
42     mask_real = out[:, 0]
43     mask_imag = out[:, 1]
44     mask_real = F.pad(mask_real, [0, 0, 1, 0])
45     mask_imag = F.pad(mask_imag, [0, 0, 1, 0])
46
47     # 极坐标系复值掩蔽的计算和滤波
48     if self.masking_mode == 'E':
49         mask_mags = (mask_real ** 2 + mask_imag ** 2) ** 0.5
50         real_phase = mask_real / (mask_mags + 1e-8)
51         imag_phase = mask_imag / (mask_mags + 1e-8)
52         mask_phase = torch.atan2(imag_phase, real_phase)
53         mask_mags = torch.tanh(mask_mags)
54         est_mags = mask_mags * spec_mags
55         est_phase = spec_phase + mask_phase
56         real = est_mags * torch.cos(est_phase)
57         imag = est_mags * torch.sin(est_phase)
58         ...
59
60         out_spec = torch.cat([real, imag], 1)
61
62         # 波形合成部分
63         out_wav = self.istft(out_spec)
64         out_wav = torch.squeeze(out_wav, 1)
65         out_wav = torch.clamp_(out_wav, -1, 1)
66     return out_spec, out_wav

```

代码通过一维卷积实现STFT和iSTFT操作，幅度和相位可以通过实部和虚部计算得到，即 $|X| = \sqrt{X_r^2 + X_i^2}$ ,  $\angle X = \arctan(\frac{X_i}{X_r})$ 。为了数值稳定， $\sqrt{(\cdot)}$ 内加了小的正值。由于DCCRN舍弃了直流量，因此在送入网络前特征的直流成分被截断，网络输出结果在直流频点用0填充。

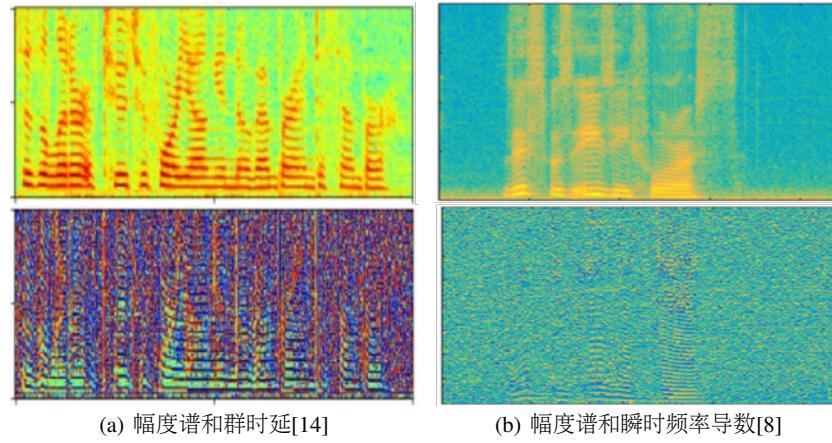
## 2.4 基于幅度相位分量估计的语音增强

由于基于幅度谱的语音增强算法的成功应用，直接估计幅度和相位成分是一个极其自然的想法。尽管前文提到的估计实部和虚部的方式理论上完成了相同任务的同时还规避了相位优化的问题，然而训练目标理论上的等价并不代表网络的优化会符合我们的预期，D. Yin等就指出Ephrat等将cIRM作为目标的算法(如图2.3-2)估计的掩蔽虚部分量几乎全为0[13]，因此仍有一些工作在相位恢复问题上继续努力。

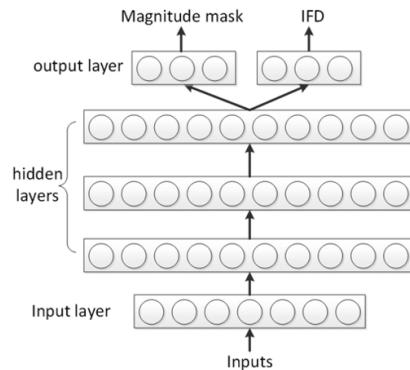
需要注意的是，目前还缺乏相关文献表明更一般的cIRM以及基于极坐标系下的复值掩蔽存在相似的问题。另外，除了训练目标外，损失函数同样对网络增强结果存在影响。只能说目前关于复数谱恢复的问题尚存在提升的空间，在某些复值掩蔽和损失函数的作用下存在和预期相距较大的情况。但更一般的结论仍需等待更多的实验结果。

2018年，一些相位估计网络被提出用于解决相位问题，如N. Takahashi等人提出将相位离散化作为一个分类问题，S. Takamichi等人则在增强的幅度谱的基础上用一个von-mises分布全连接网络在相位距离和群时延距离的损失函数下估计相位谱。

解决相位无清晰结构的另一种直观的想法就是寻找能够表征相位信息的有清晰结构模式的物理量。2019年，Z. Wang等和N. Zheng等分别提出同时估计幅度谱和群时延/瞬时频率导数，成为真正意义上基于深度学习的幅度相位分量估计算法。群时延和瞬时频率导数等物理量其实早在传统语音增强算法中就被提出用于代替相位谱的估计，通过图2.8也可以发现，两个物理量的结构与幅度谱结构也具有极强的相关性。两个工作的网络结构十分类似，如图2.9所示，将基于幅度谱的语音增强模型(LSTM/DNN)的最后一层拆分成两个分支，分别估计幅度和与相位相关的物理量。



**Fig. 2.8** 群时延和瞬时频率导数示意图



**Fig. 2.9** 基于幅度相位分量估计的语音增强模型架构[]

2020年，PHASEN按如下方式表征语音增强过程：

$$\widehat{S} = |X| \cdot \widehat{M} \cdot \Psi. \quad (2.6)$$

其采用双流结构将幅度谱掩蔽和复数谱映射的方式结合起来，一个分支用于和基于幅度谱的语音增强算法一样估计幅度掩蔽 $\widehat{M}$ ，另一个分支则只映射与相位相关的复谱 $\Psi$ ，与基于复谱映射的方法一样， $\Psi$ 被分解为实部 $\Psi_r$ 和虚部 $\Psi_i$ 两部分，通过两个分支之间的信息交互实现相位的预测。

## 2.5 小结

基于复数谱的语音增强方法进一步提升了基于幅度谱的语音增强算法性能，并成为目前语音增强研究的重点。这一阶段编解码结构的性能逐渐成为多数人的共识，许多高效的网络模块相继被提出，双流结构等新的结构也开始出现。严格地说，后续的第4章中的大部分方法仍属于基于复数谱的语音增强算法。然而由于其重框架轻结构，因此单列一章。



## Chapter 3

### 混合式和生成式语音增强模型

#### 3.1 混合式语音增强模型

神经网络除了用于直接估计滤波器系数或语谱，也常作为语音增强系统的参数估计器，我们将这类算法称为混合式语音增强(Hybrid speech enhancement)模型。

混合式语音增强模型主要建立在传统的基于统计模型的语音增强算法的基础上，利用神经网络实现先验信噪比的估计。这是由于，基于统计模型的语音增强算法的增益因子往往由先验和后验信噪比决定，而后验信噪比可以通过判决引导(Decision-Directed, DD)算法从先验信噪比中估计得到。因此，先验信噪比的估计是这类算法的核心。由于传统先验信噪比估计算法面对高度非平稳噪声缺乏良好的跟踪能力，神经网络被用于代替对应的传统模块以实现不对噪声和语音特性进行假设的、能够迅速响应噪声特性变化的先验信噪比估计。早在2011年，利用神经网络估计先验信噪比的算法就已被提出[15]用于校正传统算法估计的先验信噪比。2018年，Xia等人利用RNN辅助DD算法估计先验信噪比。2019年和2020年，DeepXi和DeepMMSE分别被提出利用网络估计先验信噪比的累计分布函数以改善网络收敛性能。

这里以DeepXi为例介绍一下基于先验信噪比估计的混合式语音增强模型的大致流程。DeepXi建立在传统的MMSE-LSA算法之上，MMSE-LSA的增益函数表示为：

$$G(t, f) = \frac{\xi(t, f)}{\xi(t, f) + 1} \exp\left(\frac{1}{2} \int_{v(t, f)}^{\infty} \frac{e^{-t}}{t} dt\right), \quad (3.1)$$

式中  $v(t, f) = \frac{\xi(t, f)}{\xi(t, f) + 1} \gamma(t, f)$ ，取决于先验信噪比  $\xi(t, f)$  和后验信噪比  $\gamma(t, f) = \xi + 1$ 。可以看到，增益因子中所有参数均取决于先验信噪比  $\xi$ 。关于算法增益函数的推导这里不做展开，感兴趣的读者可参看[16]。

DeepXi采用ResLSTM模型训练带噪幅度谱和先验信噪比的累计分布函数  $\bar{\xi} = \frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{\xi_{dB} - \mu}{\sqrt{2}\sigma}\right) \right]$  之间的映射。在推理阶段，估计的先验信噪比  $\hat{\xi}$  由网络输出  $\bar{\xi}$  经过  $\hat{\xi}_{dB} = \sqrt{2}\sigma \operatorname{erf}^{-1}(2\hat{\xi} - 1) + \mu$  及  $\hat{\xi} = 10^{(\hat{\xi}_{dB}/10)}$  得到后，

代入MMSE-LSA增益函数公式3.1即可对带噪幅度谱进行滤波得到增强语谱 $\hat{S} = G \cdot |X| \cdot e^{\angle X}$ 。

另一类混合式语音增强算法严格意义上并不仅在幅度谱上进行，也并非源于语音增强领域的传统算法。RNNoise是其中最著名的代表之一，其改进版PercepNet更是取得了令人瞩目的性能。他们明显受到了语音编码算法的影响，按照ELT编码器的思路，将语谱分解为语谱包络和语谱细节两部分。根据语音的生成机制，语音的包络由声道形状的变化得到，语谱细节则是由声带的准周期性震动产生，主要是基频及其谐波。根据人耳听觉感知机制，语谱包络可以在某种听觉感知相关的低分辨率频带(band)上进行操作。CELT中“应当保证重构的语谱包络能量在频带上与纯净语谱包络相等”的思想启发我们可以通过计算增益函数得到增强语音的谱包络，此时的谱包络在浊音段仍是粗糙的，而在其他段则与纯净语谱听感上相似。而语谱细节——也就是基频及其谐波——可以通过基频滤波进行修正，抑制掉谐波间的噪声，此时应在高分辨率的STFT谱或时域上进行，主要包含两步：首先利用梳状滤波器得到滤波后的语谱 $\hat{P}$ ，然后根据网络估计的基频滤波强度 $r$ 加权组合带噪谱成分和滤波信号谱成分以增加语音的自然度 $Z = (1 - r) \cdot X + r \cdot \hat{P}$ 。图3.1展示了上述过程，其中特征提取器feature extraction、神经网络DNN model和包络后滤波envelope postfilter承担了语谱包络估计任务(分别用于变换到低分辨率的ERB频带、估计谱包络增益函数和通过增益补偿和加混响改善听感)；而神经网络DNN model和基频滤波pitch filtering则承担了语谱细节估计任务(分别用于估计基频滤波强度和基频滤波)。毫无疑问神经网络并非混合式语音增强模型的重点，因此这里一笔掠过：RNNoise试图按照谱减法框架设计RNN网络完成谱包络修正，PercepNet则直接使用多任务学习的CRN同时估计谱包络增益因子和基频滤波强度系数。

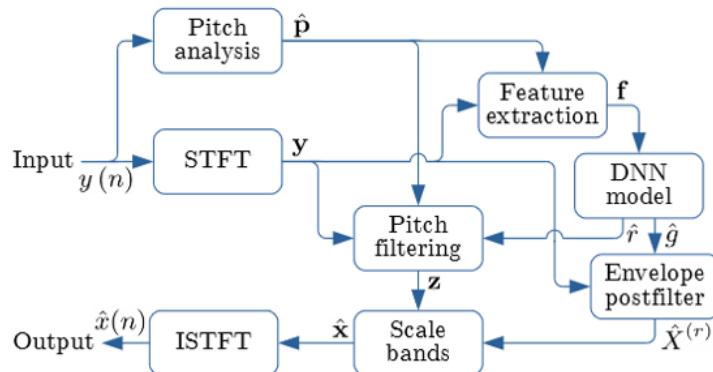


Fig. 3.1 PercepNet框图[?]

值得注意的是，关于RNNoise和PercepNet的讨论并未结束，其将语音分解为包络和细节的思想被后续工作完全神经网络化，这种思想催生出的模型在

第4章中详细介绍。此外，随着超宽带和全带语音增强的任务得到重视，其在低分辨率听觉感知频带的处理思路仍具参考价值。

除了以上两类算法之外，还有一种时域混合式语音增强算法。这类算法将神经网络与Kalman滤波的结合，将语音生成视为一种自回归(auto-aggressive, AR)过程(这同样在语音编码和生成领域被广泛应用)，使用基于MMSE准则的Kalman滤波器增强语音波形，从而隐式地修正带噪语音的幅度和相位。

这一过程中，Kalman滤波器的参数主要是线性预测系数(linear prediction coefficients, LPCs)以及激励噪声和测量噪声的统计特性。神经网络被用于估计Kalman滤波器的三个参数(也可能只估计LPCs，其他参数利用传统算法估计)，这些参数被带入Kalman滤波器的“预测-更新”公式，从而得到滤波器系数。语音增强任务的“预测-更新”公式如下，同样，更详细的公式化过程可参看[17]：

$$\begin{aligned}\mathbf{K}(n) &= \mathbf{P}(n|n-1)(\mathbf{R}_w + \mathbf{P}(n|n-1))^{-1}, \\ \hat{\mathbf{s}}(n|n) &= \hat{\mathbf{s}}(n|n-1) + \mathbf{K}(n)(\mathbf{x}(n) - \hat{\mathbf{s}}(n|n-1)), \\ \mathbf{P}(n|n) &= (\mathbf{I} - \mathbf{K}(n))\mathbf{P}(n|n-1), \\ \hat{\mathbf{u}}(n+1|n) &= \mathbf{F}\hat{\mathbf{u}}(n|n), \\ \mathbf{P}(n+1|n) &= \mathbf{F}\mathbf{P}(n|n)\mathbf{F}^T + \sigma_v^2 \mathbf{G}\mathbf{G}^T,\end{aligned}\quad (3.2)$$

式中  $\mathbf{s}(n) = [s(n-p+1), \dots, s(n)]^T$ ,  $\mathbf{x}(n) = [x(n-p+1), \dots, x(n)]^T$ ,  $\mathbf{G} = [0, \dots, 0, 1]^T$ ,  $p$  是对纯净语音进行AR建模的阶数。 $\mathbf{F}$ 、 $\sigma_v^2$  和  $\mathbf{R}_w$  分别代表包含LPCs的变换矩阵、AR建模时的激励噪声方差和kalman滤波建模时的测量噪声协方差矩阵。

通过观察  $\mathbf{s}(n)$  的定义不难得出，采样点  $n$  的增强语音  $\hat{\mathbf{s}}(n)$  可通过  $\mathbf{G}\hat{\mathbf{s}}(n|n)$  得到。

上面介绍了最简单的基于Kalman滤波的混合式语音增强模型，利用线谱频率(line spectrum frequencies, LSFs)代替LPCs作为网络映射目标、采用子带Kalman滤波和Colored-Noise Kalman滤波等Kalman滤波器变体等改进也相继被提出。图3.2展示了这类算法的训练和推理流程。

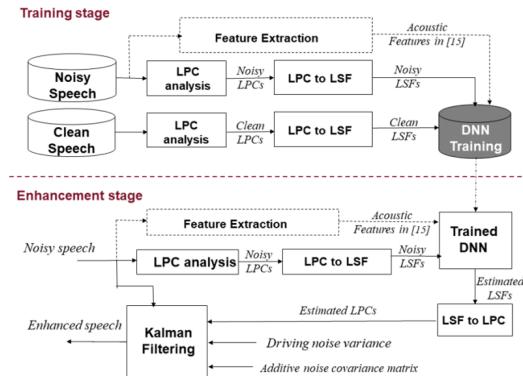


Fig. 3.2 基于Kalman滤波的混合式语音增强框图[?]

### 3.2 生成式语音增强模型

1. GAN
2. VAE
3. Diffusion Model
4. TTS-Style SE

### 3.3 小结

## **Chapter 4**

### 基于解耦的语音增强



## **Chapter 5**

### 基于时域的语音增强



## Chapter 6

### 语音增强算法的因果性

语音增强技术根据其应用的不同可以分为实时系统和非实时系统两大类。助听辅听、耳机手机通话等是典型的实时应用场景，其中的语音增强算法必须在使用较少的未来信息甚至不使用未来信息的情况下(即模型是因果的)下完成处理，同时，算法在相应的硬件设备上的计算时间也应满足实时要求。需要注意的是，这里的实时是一个相对的概念，绝大多数情况的语音增强并不要求逐采样点的实时处理，而是指在较短延时内完成处理，而对于不同的应用，能容忍的延时也各不相同。由于不同硬件平台的算力不同，更常见的是用模型的因果性代替对其实时性的评价，并使用参数量和计算量作为是否能够在某个硬件上实时处理的参考指标。这是由于语音增强的应用场景是如此广泛、模型压缩技术仍有对模型效率改进的可能、以及芯片等计算资源的迅猛发展，因此即使某个提出的语音增强模型无法应用在某个平台上，该模型可能仍不失其价值。

语音增强处理延时的定义并不止一种，这里给出一种常见的定义：对于采用重叠相加(overlap-add)方法重构的频域语音增强算法，时延由OLA带来的时延(合成窗长)、频域变换带来的时延(帧移)以及算法接收的未来信息带来的时延三部分组成，另外算法的处理时间应不大于帧移。

为了保证模型比较的公平性以确定模型中每个模块的有效性，模型因果与否是应该明确的。然而，很遗憾的是，社区内常出现工作有意无意地忽视了这一问题，从而误导了社区对于某种特征/模块/损失函数的认识。因此，作为一本语音增强导读，有必要将这部分单列一节进行说明。

这一话题无疑是必要但敏感的，因此这里不得不表明个人态度。其一，本节的立意是尽量澄清一些社区内的误区、尽量减少初学者所走的弯路、尽量提升初学者对文章和开源代码的判断能力，而非否定以往的和未来的任何工作。其二，尽管一些工作并不严格满足因果处理，但与因果处理保持一致的baselines比较时大约是能够保证其创新部

分的有效性的。另一方面，尽管个人由于不希望看到语音前端陷入到和计算机视觉一样门槛大幅度降低、滥竽充数者比例渐高、工业界多方面压制学术界的状况，因此对开源并不持支持态度。但无疑开源者是对社区发展有着推动和帮助作用的，他们的贡献是巨大的。窃以为不宜以这种很可能无意的失误而使贡献者受到责难。我也反对在任何高速发展的学科因为小的技术失误而全面否定整个工作的行为，这对学科和个人的发展是有负面作用的(很不幸近年来这种趋势似乎是明显的)。更何况实时语音增强中大量工作并不因这一问题而失去其价值，而且语音增强在非实时场景同样也有广泛应用。

## 6.1 语音音量的归一化

由于说话人本身音量以及说话人距传声器距离远近的差异，语音音量大小的动态范围通常很大。为了保证语音增强模型能够处理不同音量的语音，对训练数据的处理通常是重要的。社区内对此常见的两种非因果操作是对训练和测试过程中的带噪语音和纯净语音进行能量均方根归一化(即假设语音经过了理想的自动增益控制(AGC))、以及利用整段带噪音频计算均值和方差后用作均值方差归一化(这种方式既有利用训练集的均值方差对测试语料做归一化的，又有直接测试语料计算均值和方差做归一化的，显然后者是非因果的)。以这种方式压缩了数据的动态范围，方便网络优化。替代这类非因果操作的方式有两种：一种是对输入特征进行在线归一化、另一种是在损失端进行归一化。

前者常见的有在线频率依赖的均值方差归一化，当前帧的均值方差均以指数衰减平滑的形式估计得到：

$$\begin{aligned}\mu(t, f) &= \alpha\mu(t-1, f) + (1-\alpha)|X(t, f)|, \\ \sigma^2(t, f) &= \alpha\sigma^2(t-1, f) + (1-\alpha)|X(t, f)|^2, \\ |X^{norm}(t, f)| &= \frac{|X(t, f)| - \mu(t, f)}{\sqrt{\sigma^2(t, f) - \mu^2(t, f)}}.\end{aligned}\tag{6.1}$$

上述关系式代码为：

```

1 # https://github.com/GuillaumeVW/NSNet/blob/master/
2     dataloader/wav_dataset.py lines 85 - 102
3 # mean normalization
4 frames = []
5 x_lps = x_lps.transpose(0, 1)
6 n_init_frames = self.n_init_frames
7 alpha_feat_init = self.alpha_feat_init
8 alpha_feat = self.alpha_feat
9 for frame_counter, frame_feature in enumerate(x_lps):
10 if frame_counter < n_init_frames:
11     alpha = alpha_feat_init
12 else:
```

```

12     alpha = alpha_feat
13     if frame_counter == 0:
14         mu = frame_feature
15         sigmasquare = frame_feature.pow(2)
16         mu = alpha * mu + (1 - alpha) * frame_feature
17         sigmasquare = alpha * sigmasquare + (1 - alpha) *
18             frame_feature.pow(2)
19         sigma = torch.sqrt(torch.clamp(sigmasquare - mu.pow(2),
20             min=1e-12)) # limit for sqrt
21         norm_feature = (frame_feature - mu) / sigma
22         frames.append(norm_feature)

```

此外，其他累积归一化方式和可学习归一化方式也被探索。

后者则是将语音能量均方根归一化的操作转移到损失端，整句语音的能量仅在训练阶段被使用，因此模型是因果的。这类方法被用于复数域语音增强模型，其损失函数被定义为[7]：

$$\mathcal{L} = \frac{1}{\sigma} (\alpha \sum_{t,f} |S - \hat{S}|^2 + (1 - \alpha) \sum_{t,f} ||S| - |\hat{S}||^2), \quad (6.2)$$

其中 $\sigma$ 是指经过语音活动检测(VAD)后的语音能量。一般同时还会在训练阶段对语料音量以某种随机分布随机缩放，如在训练阶段将纯净语音和噪声的音量都服从均值为-26 dBFS，方差为10 dB的高斯分布进行随机缩放[7]，类似操作的代码见于DPCRN3：

```

# data_loader.py lines 163 - 192
1 gain = np.random.normal(loc=-5, scale=10)
2 gain = 10**((gain/10)
3 gain = min(gain, 3)
4 gain = max(gain, 0.01)
5 ...
6 batch_clean[i, :] = clean_s * gain
7 batch_noisy[i, :] = noisy_s * gain

```

## 6.2 因果模块讨论

网络层面的因果性相对更易察觉，不过这有时要求对已有的深度学习API有足够的认知。目前基于深度学习的语音增强模型常用层等。本节以Pytorch框架为例简要地对这些模块进行讨论。这里默认读者已熟悉这些模块并熟悉相关的Pytorch函数，使用其他框架的读者可做参考，结合所使用框架的文档进行核对。

1. LSTM层和GRU层：torch.nn.LSTM和torch.nn.gru中的参数bidirectional设为False(默认)即只使用历史信息；
2. 二维卷积层：需对torch.nn.Conv2d/Conv1d输入特征的时间帧维度向历史帧方向补零(时间帧方向卷积核不为1时)；

3. 二维转置卷积：需对`torch.nn.Transpose2d`输出结果的时间帧维度进行在未来帧方向截断(时间帧方向卷积核不为1时);
4. 注意力机制模块：需对点积相似度公式中的`Softmax`函数内加掩蔽项，这样虽然是因果的，但在推理时内存和运行时间仍会随帧数增加而增长。基于`chunk`或类`TransformerXL`的`attention`才有可能满足实时处理要求。
5. 标准化层：`Batch Normalization`应在推理时将模型设为`eval()`模式；`Instance Normalization`应将参数`track_running_stats`设置为`True`。
6. 池化层：因果模型中不应对时间帧维度进行池化。