

内蒙古大学计算机学院“语音理解与生成”研究组

硕士研究生招生【长期有效】

内蒙古大学计算学院“语音理解与生成”研究组（Speech Understanding and Speech Generation Research Group, S2Group）在刘瑞研究员的带领下长期从事深度学习、人工智能、语音信息处理、表现力语音合成等相关工作。团队成员拥有雄厚理论研究积累，相关成果发表于 SCI 一区 Top 期刊 IEEE/ACM Transactions on Audio, Speech, and Language Processing (IEEE/ACM TASLP)，语音领域顶级会议 ICASSP、INTERSPEECH 等。

团队介绍

刘瑞（团队负责人）

- 内蒙古大学研究员、博士生导师
- 新加坡国立大学联合培养博士，新加坡国立大学博士后
- IEEE/ACM-TASLP、IEEE Signal Processing Letters、ICASSP、INTERSPEECH、Blizzard Challenge 等多个领域重要期刊和会议审稿人
- O-COCOSDA 2021, IWSDS 2021, SIGDIAL 2021 等学术会议组织主席
- 2022 年全国人机语音通讯学术会议 (NCMMSC2022) 工业联络主席
- 中国计算机学会语音对话与听觉专委会委员、中国人工智能学会青年工作委员会委员
- 电气和电子工程师协会 (IEEE)、国际语音通讯学会 (ISCA)、国际计算机学会 (ACM)、中国计算机学会 (CCF)、中国人工智能学会 (CAAI) 会员

团队依托内蒙古自治区蒙古文信息处理重点实验室（主任：飞龙教授）和蒙古文智能信息处理技术国家地方联合工程研究中心（主任：高光来教授）开展研究，S2Group 研究组拥有 2 名博士研究生及 10 名硕士研究生，已经逐步形成了一个稳定的研究梯队。

研究成果

项目：团队目前承担“2022 内蒙古大学高层次人才引进项目”以及“2022 国家自然科学基金青年科学基金项目”，另外参与多项国家自然科学基金面上项目、国家重点研发计划项目、国家自然科学基金地区科学基金项目、新加坡国防科技部重点项目等。

论文：团队在国内外人工智能及语音信息处理领域顶级期刊和会议上发表论文 30 余篇。包括 5 篇 SCI 一区 Top 期刊（4 篇 IEEE/ACM TASLP 和 1 篇 IEEE Internet of Things Journal）和 2 篇 SCI 二区期刊（Neural Networks 和 IEEE Signal Processing Letters），以及若干篇 ICASSP、InterSpeech 会议。论文累计引用 300 多次（Google Scholar, H-index=11），引用者包括来自美国卡耐基梅隆大学、英国剑桥大学、英国爱丁堡大学、日本名古屋工业大学、新加坡国立大学、新加坡科技与设计大学、中科院自动化所、香港中文大学、清华大学、西北工业大学等研究机构的国内外知名学者。

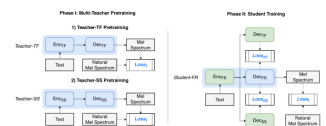


Fig. 1. Architecture of a multi-task knowledge distillation (MTKD) scheme in a shared TTS system. The parameters in this box have an identical value across tasks, while those in gray boxes are trained by using the pre-trained parameters from those T

B. Phase II: Student Training

At the same time, the student decoder $D_{S, \text{pub}}$ takes the previous hidden states h_{t-1} , estimated speech \hat{y}_{t-1} , and the attention score α_{t-1} as input, and predicts the hidden state h_t at each step t .

$$h_t = D_{S, \text{pub}}(h_{t-1}, \hat{y}_{t-1}, \alpha_{t-1}) \quad (12)$$

Finally, we follow (8) to generate the output speech \hat{y}_t .

The training is supervised by these three objective functions:

$$Loss_{\text{pub}}, Loss_{\text{priv}}, \text{ and } Loss_{\text{TD}}$$

The two distillation loss $Loss_{\text{pub}}$ and $Loss_{\text{priv}}$ ensure that the hidden states of $D_{S, \text{pub}}$ are as close to $D_{S, \text{pub}}$ and $D_{S, \text{priv}}$ as possible. We apply two classification loss functions $Loss_{\text{pub}}$ and $Loss_{\text{priv}}$ to ensure that the hidden states of $D_{S, \text{pub}}$ are as close to those of the teacher models. At the same time, we adopt the teacher loss $Loss_{\text{TD}}$ to ensure that the predicted speech is as close to the reference natural speech. The decoder is trained in a self-supervised manner by comparing the predicted speech during training with the reference natural speech. We also formulate the training of $D_{S, \text{pub}}$ as follows:

- The trained decoder $D_{S, \text{pub}}$ takes the previous hidden state h_{t-1} , estimated speech \hat{y}_{t-1} , and the attention score α_{t-1} as input, and outputs the hidden state h_t at each time step t of the input.
- Loss: $Loss_{\text{pub}} = D_{S, \text{pub}}(h_{t-1}, \hat{y}_{t-1}, \alpha_{t-1})$

At the same time, the student decoder $D_{S, \text{priv}}$ takes the previous hidden states h_{t-1} , estimated speech \hat{y}_{t-1} , and the attention score α_{t-1} as input, and outputs the hidden state h_t at each time step t of the input.

Loss: $Loss_{\text{priv}} = D_{S, \text{priv}}(h_{t-1}, \hat{y}_{t-1}, \alpha_{t-1})$

At the same time, the student decoder $D_{S, \text{TD}}$ takes the previous hidden states h_{t-1} , estimated speech \hat{y}_{t-1} , and the attention score α_{t-1} as input, and outputs the hidden state h_t at each time step t of the input.

Loss: $Loss_{\text{TD}} = D_{S, \text{TD}}(h_{t-1}, \hat{y}_{t-1}, \alpha_{t-1})$

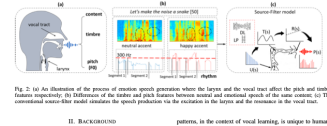


Fig. 2. An illustration of the process of emotion speech generation where the latent and the vocal tract affect the pitch and timbre features respectively. (a) Difference of the latent and pitch features between natural and emotional speech of the same content. (b) The conventional vocoder that would generate the speech parameters by the influence of the latent and the response in the vocal tract.

III. BACKGROUND

We provide here a brief primer on the conventional vocoder. The process is a simplified illustration of the process of speech generation. The speech excitation comes from the vibration of vocal folds in the larynx. The generated vocal tract is then modified by resonance of the vocal tract (nasal, oral and nasal cavities). The speech signal contains three main information components: language content, timbre, pitch and dynamics. The emotional features are embedded in these components in different ways.

Content being in the language model in speech research. The phenomena in the basic unit of speech content is most important. In this sense, there is a particular feature content. This different generation appear in different steps in the conventional VCT, that is, the segmentation of the content being in different content level (nasal, oral and nasal cavities). This content filter model needs retain the content information while converting the emotional features.

Timbre is reflected by the formant, which is a peak of the spectral envelope that results from an acoustic resonance of the human vocal tract. The timbre can represent the tone color or unique quality of a sound which can be measured objectively and classified using various, such as individual people or natural languages. In the conventional vocoder, the high-amplitude and higher-order formants are usually generated from the low-amplitude and lower-order formants. As shown in Fig. 2(b), the segmentation of the content being in different content level (nasal, oral and nasal cavities) shows that the happy voice has a larger segmentation in some words, a higher frequency response in the oral cavity, and a larger amplitude than the natural one. This is necessary to extract the emotional information from the input features.

Pitch is an important parameter in emotional speech processing systems. Emotional pitch is generated by the larynx and modulated primarily by fine changes in the tension of the vocal folds. The ability to manipulate and flexibly control pitch is essential in the content of vocal speech. It is more important in the content of vocal speech, which is a key factor in the content of vocal speech. It is more important in the content of vocal speech, which is a key factor in the content of vocal speech.



Fig. 3. Block diagram of the proposed training strategy. Teacher T_{pub} : A speech emotion recognition (SER) model is trained separately to serve as an auxiliary model to extract deep style features. A critic architecture W_{crit} is compared between the deep style features of the generated and reference speech to

speech primarily is not straightforward. One of the ways to describe a prosodic style is to show its envelope. The use of such values (1) shows a way to compare two prosodic styles (2) and (3). The envelope of a prosodic style is a measure of the amplitude of the speech signal over time. The envelope of a prosodic style is a measure of the amplitude of the speech signal over time. The envelope of a prosodic style is a measure of the amplitude of the speech signal over time.

Prosodic representation in emotional speech (PS) models. Prosodic representation in emotional speech (PS) models is a key factor in the content of vocal speech. It is more important in the content of vocal speech, which is a key factor in the content of vocal speech.

Prosodic representation in emotional speech (PS) models. Prosodic representation in emotional speech (PS) models is a key factor in the content of vocal speech. It is more important in the content of vocal speech, which is a key factor in the content of vocal speech.

系统：团队基于最新的非自回归语音合成技术搭建了蒙古语语音合成系统 MonTTS，可以用于深度学习相关研究及工业系统的开发。



图 2. MonTTS 整体框架图
 包括 (a) 模型结构；(b) 音素级声学调节器内部结构及相应的损失函数；(c) 基于蒙古语语音识别的对齐方法；(d) 基于蒙古语自回归语音合成的对齐方法。

竞赛：团队组织或参与多个国内外语音合成竞赛任务。举办“NCMMSC2022 面向蒙古语的低资源语音合成竞赛”，在国际语音合成竞赛 Billzard Challenge 2019/2020 中取得优异成绩。

获奖：在 IALP2021 学术会议发表的论文荣获会议最佳论文奖 (Best Paper Award)。



学术交流：团队成员赴美国、印度、新加坡等地参加国际学术会议，另外，受中国计算机学会、中国人工智能学会等邀请多次进行学术报告。

语音合成专题学术论坛

主办单位：语音及语言信息处理国家重点实验室
 中国计算机学会语音识别专委会

活动时间：2021年12月4日周六9:00-11:45
 参会方式：腾讯会议（会议号974378789）
 或扫描右侧二维码

会议安排：
 09:00-09:05 活动致辞（CCF语音识别与听觉专委会副主任 俞航）
 09:05-09:45 刘瑾：端到端语音合成中的韵律、情感建模研究
 09:45-10:25 胡亚东：语音合成中的韵律表征解耦和建模
 10:25-11:05 郑旭：达到语音水平的文本到语音合成
 11:05-11:45 王鑫：从语音安全的两个问题看语音合成

Speech home

语音之家公开课

语音合成中的情感强度建模研究

蒙古语智能信息处理

技术及其应用

招生说明

招生要求：

- **针对推免生，应获得所在学校当年推免资格**
- 身体健康，心理健康状况良好，思想积极乐观
- 数理基础优良，Coding 能力和计算机基本专业课（包括但不限于算法、数据结构、操作系统）优良，逻辑思维清晰
- 具有较强的外语听、说、读、写能力
- 具有较强的调查研究、综合分析问题、解决问题的能力
- **优先条件**：英语要求通过全国大学英语四级考试（成绩 425 分以上）
- **优先条件**：具有搭建/训练神经网络模型的经验，熟悉 Tensorflow、Pytorch 等深度学习框架。
- **优先条件**：在领域重要会议上发表过一篇或以上学术论文（个人排名前三）

学习待遇：

- 指导：每名同学受团队负责人直接指导，并配备一名高年级研究生进行指导
- 设备：提供研究所需必要硬件设备 (高性能 GPU 服务器)
- 场地：提供学习座位，同时参与实验室学术讨论等活动
- 福利：丰厚的奖学金制度
- 交流：支持参加国内外学术会议

交流合作：实验室与多家海内外大学和科研机构建立了长期稳定的学术及项目合作关系，如新加坡国立大学、英国帝国理工学院、香港中文大学、华南理工大学、腾讯公司、微软公司等。可选派优秀学生到上述研究机构以及国际公司访问学习或科研项目合作。

报名方式

请发送电子版简历至负责人 刘瑞 研究员（邮箱：imucslr@imu.edu.cn 或 liurui_imu@163.com）

更多团队相关信息请访问团队官网查看~

https://ttslr.github.io/index_S2Group.html

相关信息

刘瑞个人主页：



S2Group 主页：



“智能语音新青年”公众号



同时欢迎 **在读本科生** 加入研究组进行科研实践！