

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】

内蒙古大学计算机学院“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。

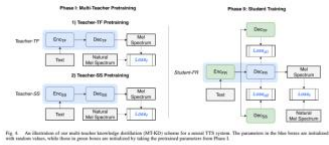


Fig. 4. Architecture of the multi-task knowledge distillation (MTKD) framework for a neural TTS system. The parameters in the blue boxes are inherited from the teacher model, while those in green boxes are initialized by using the pre-trained parameters from the teacher.

B. Phase II: Student Training

The Student II model has the same network architecture as the teacher model, which includes E_{T2} and D_{T2} . On the right part of Fig. 4, we illustrate the schematic diagram of the training process. To transfer from the pre-trained models, we use E_{T2} or D_{T2} to initialize E_{S2} , which generates the linguistic encoding from input text. In a pilot experiment, we found an difference between the two initialization methods in some of utterances. Therefore, we adopt E_{T2} to initialize the student model E_{S2} in future experiments.

Next, the student decoder D_{S2} takes the given input text $x = [x_1, x_2, \dots, x_N]$ and outputs linguistic encoding $h = [h_1, h_2, \dots, h_N]$:

$$h = E_{S2}(x) \quad (6)$$

We train the Student II model following the *Teacher2Student* process. During training, the teacher decoder and student decoder take the same encoder output sequence as the input, and generate their own hidden states sequentially. We use two distillation loss functions, L_{enc} and L_{dec} , to supervise the hidden states generated by the student decoder to be close to those of the teacher models. At the same time, we adopt the Student loss function L_{S2} to ensure that the predicted speech is close to the reference natural speech. The distillation loss is a loss measuring model that is consistent with the distillation process during our time inference. We now formulate the training of D_{S2} .

The pre-trained D_{T2} takes the previous hidden state h_{t-1} , current speech frame s_t , and the attention score α_t as input, and outputs the hidden state h_t at each time step t by (7):

$$h_t = D_{T2}(h_{t-1}, s_t, \alpha_t) \quad (7)$$

The Student decoder D_{S2} receives another hidden state h_{t-1} at each time step t by following (8):

$$h_t = D_{S2}(h_{t-1}, s_t, \alpha_t) \quad (8)$$

At each time step, only Student II model is involved, where we use E_{T2} to process input text, and D_{S2} to generate an acoustic:

$$s_t = D_{S2}(E_{T2}(x), s_t, \alpha_t) \quad (9)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (10)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (11)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (12)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (13)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (14)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (15)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (16)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (17)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (18)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (19)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (20)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (21)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (22)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (23)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (24)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (25)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (26)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (27)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (28)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (29)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (30)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (31)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (32)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (33)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (34)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (35)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (36)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (37)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (38)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (39)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (40)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (41)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (42)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (43)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (44)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (45)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (46)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (47)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (48)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (49)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (50)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (51)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (52)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (53)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (54)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (55)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (56)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (57)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (58)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (59)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (60)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (61)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (62)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (63)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (64)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (65)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (66)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (67)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (68)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (69)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (70)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (71)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (72)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (73)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (74)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (75)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (76)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (77)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (78)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (79)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (80)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (81)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (82)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (83)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (84)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (85)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (86)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (87)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (88)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (89)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (90)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (91)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (92)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (93)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (94)$$

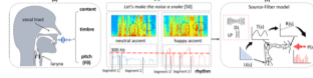


Fig. 2. (a) Illustration of the process of emotion speech generation where the layers and the vocal tract affect the pitch and timbre. (b) Illustration of the process of emotion speech generation where the layers and the vocal tract affect the pitch and timbre. (c) Illustration of the process of emotion speech generation where the layers and the vocal tract affect the pitch and timbre.

We provide here a brief primer on the emotional speech, the speaker's emotion and the emotional voice-acted speech theory.

B. BACKGROUND

A. Adaptation in Emotional Speech

Fig. 2 (a) gives a qualitative discussion of the process of speech generation. The speech is generated from the vibration of vocal cords in the larynx. The generated vocal sound is then modulated by resonance of the vocal tract (pharynx, oral cavity, nasal cavity). The speech signal contains four main information components: language content, intonation, pitch, and timbre. The emotional features are embedded in each component in different way.

Content belongs to the language model in speech research. The phoneme is the basic unit of speech content in most languages. Each phoneme has a particular formant pattern. This different phoneme appear at different places in the spectrum. As shown in Fig. 2 (b), the spectrum of the same content speech in different phrases have similar shapes to be kept unchanged. Thus, the voice filter model needs to retain the content information while converting the emotional features.

Timbre is referred by the formant, which is the peak of the spectral envelope that results from an acoustic resonance of the vocal tract. The timbre is related to the timbre of the sound. In the emotional speech, the high amount of a happy or angry voice tends to vocal chords and higher than the low amount or neutral voice [36]. As shown in Fig. 2 (c), the comparison between a happy and a neutral utterance shows that the happy voice has a deeper spectrum in some words, and a higher formant frequency range in the rest of the utterance.

Which is a useful feature between the two different loss terms. Algorithm 1 describes the complete training process of the proposed MTKD.

Both a self-supervised and a supervised learning framework are used to train the model. The self-supervised learning is used to train the model, while the supervised learning is used to train the model.

At each time step, only Student II model is involved, where we use E_{T2} to process input text, and D_{S2} to generate an acoustic:

$$s_t = D_{S2}(E_{T2}(x), s_t, \alpha_t) \quad (9)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (10)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (11)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (12)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (13)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (14)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (15)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (16)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (17)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (18)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (19)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (20)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (21)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (22)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (23)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (24)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (25)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (26)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (27)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (28)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (29)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (30)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (31)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (32)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (33)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (34)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (35)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (36)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (37)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (38)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (39)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (40)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (41)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (42)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (43)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (44)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (45)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (46)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (47)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (48)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (49)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (50)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (51)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (52)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (53)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (54)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (55)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (56)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (57)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (58)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (59)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (60)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (61)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (62)$$

$$s_t = D_{S2}(E_{S2}(x), s_t, \alpha_t) \quad (63)$$

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。

- **交流合作案例 1:** 博士一年级研究生左昊麟与[启元实验室（北京）](#)赵金明研究员和[香港中文大学（深圳）](#)李海洲教授合作在 ICASSP2023 发表论文
- **交流合作案例 2:** 硕士一年级学生 张锦华 与 [香港中文大学（深圳）](#) 李海洲教授、武执正教授合作
- **交流合作案例 2:** 硕士二年级学生 胡一帆与 [香港中文大学（深圳）](#) 李海洲教授合作

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！