



دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

پایان نامه کارشناسی
مهندسی کامپیوتر

پیاده سازی یک سرویس تبدیل صدا مبتنی بر یادگیری عمیق

نگارش

کسرا دماوندی اصلی

استاد راهنما

حسین صامتی

تیر ۱۴۰۲

به نام خدا
دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

پایان نامه کارشناسی

این پایان نامه به عنوان تحقق بخشی از شرایط دریافت درجه کارشناسی است.

عنوان: پیاده سازی یک سرویس تبدیل صدا مبتنی بر یادگیری عمیق

نگارش: کسرا دماوندی اصلی

کمیته ممتحنین

استاد راهنما: حسین صامتی

امضاء:

استاد مشاور: استاد مشاور

امضاء:

استاد مدعو: استاد ممتحن

امضاء:

تاریخ:

سپاس

از استاد بزرگوارم جناب دکتر صامتی که با کمک‌ها و راهنمایی‌های بی‌دریغشان، مرا در به سرانجام رساندن این پایان‌نامه یاری داده‌اند، تشکر و قدردانی می‌کنم. همچنین از آقای سروش گوران که در تمام فرایند انجام این پروژه مرا مشاوره دادند و سایر عزیزانی که مرا یاری نمودند صمیمانه سپاسگزارم.

چکیده

تبدیل صدا فرایند تبدیل صدای یک گوینده به گوینده دیگر است در حالی که محتوای زبانی گفتار باید حفظ شود. از این فناوری در حوزه‌هایی چون رسانه، تقویت ارتباط گفتاری و سرگرمی استفاده می‌شود. امروزه غالباً پروژه‌های تبدیل صدا مبتنی بر شبکه‌های عصبی می‌باشند. در این پروژه، هدف ما پیاده‌سازی یک سرویس تبدیل صدا می‌باشد که قابلیت دریافت فایل صوتی از کاربر و ارسال نتیجه به او را دارد. Backend این سرویس با Django پیاده‌سازی شده است. همچنین با استفاده از MongoDB یک پایگاه داده برای ذخیره اطلاعات مورد نیاز ساخته شده است. در نهایت با ایجاد یک بات تلگرام این سرویس را به بهره‌برداری رساندیم.

کلیدواژه‌ها: تبدیل صدا، شبکه‌های عصبی، پایگاه داده، Backend

فهرست مطالب

۱	مقدمه	۱
۱	۱-۱ تعریف مسئله	۱
۱	۲-۱ کاربردها	۱
۲	۳-۱ انواع تبدیل صدا	۲
۲	۴-۱ اهداف پژوهش	۲
۲	۵-۱ ساختار پایان نامه	۲
۳	۲ کارهای پیشین	۳
۳	۱-۲ HiFi GAN	۳
۳	۱-۱-۲ Generator	۳
۴	۲-۱-۲ Discriminator	۴
۶	۳-۱-۲ خطای یادگیری	۶
۷	۲-۲ WavLM	۷
۹	۳ مدل FreeVC	۹
۹	۱-۳ مدل FreeVC	۹
۱۰	۲-۳ ساختار مدل FreeVC	۱۰
۱۰	۱-۲-۳ انکودر پیشین	۱۰
۱۱	۲-۲-۳ انکودر گوینده	۱۱

۱۱	۳-۲-۳ انکودر پسین
۱۲	۳-۲-۴ دکودر و جداکننده
۱۲	۳-۲-۵ داده افزایی
۱۳	۳-۳ خطای مدل
۱۵	۴ گزارش‌های فنی آموزش و استنتاج مدل
۱۵	۴-۱ تنظیمات مدل
۱۶	۴-۲ آموزش مدل
۱۸	۴-۳ استنتاج
۱۹	۵ پیاده‌سازی سرویس تبدیل صدا
۱۹	۵-۱ پایگاه داده
۲۰	۵-۲ Backend
۲۲	۶ جمع‌بندی
۲۳	آ مطالب تکمیلی

فهرست شکل‌ها

۴	۱-۲ ساختار Generator در HiFi GAN
۵	۲-۲ ساختار Discriminator ها در HiFi GAN
۷	۳-۲ ساختار WavLM
۱۱	۱-۳ سمت راست: استنتاج، سمت چپ: یادگیری
۱۳	۲-۳ داده‌افزایی مبتنی بر تغییر اندازه اسپکتروگرام
۱۶	۱-۴ یادگیری با انکودر از پیش آموزش دیده
۱۶	۲-۴ یادگیری بدون داده‌افزایی
۱۷	۳-۴ یادگیری با داده‌افزایی
۱۷	۴-۴ وارد کردن آدرس فایل‌های صوتی
۱۷	۵-۴ فایل متنی
۱۸	۶-۴ فایل convert.txt
۱۸	۷-۴ استنتاج
۱۸	۸-۴ استنتاج
۲۰	۱-۵ مجموعه حساب کاربر

فصل ۱

مقدمه

۱-۱ تعریف مسئله

امروزه هوش مصنوعی به یکی از حوزه‌های پرطرفدار تکنولوژی بدل شده است که به مسائل گوناگونی از جمله پردازش گفتار، متن، تصویر و... می‌پردازد. یکی از مسائلی که در این حوزه به آن پرداخته می‌شود تبدیل صدا است. تبدیل صدا فرایند تغییر صدای یک گفتار به صدای شخصی دیگر (گوینده مقصد) است به گونه‌ای که طبیعی بودن آن و همچنین محتوای زبانی مانند کلمات، لحن، ریتم و تأکیدها حفظ شود. این فرایند معمولاً در ۳ مرحله صورت می‌گیرد:

۱. استخراج اطلاعات مربوط به محتوای گفتار

۲. استخراج اطلاعات مربوط به صدای گوینده مقصد

۳. ساخت موج صوتی جدید با استفاده از اطلاعات به دست آمده.

۲-۱ کاربردها

از تبدیل صدا می‌توان در رسانه‌ها استفاده کرد. برای مثال در صنعت دوبله فیلم و سریال می‌توان از صدای صداپیشگان قدیمی که دیگر فعالیت نمی‌کنند استفاده کرد. از دیگر کاربردهای تبدیل صدا می‌توان به تقویت ارتباط گفتاری اشاره نمود که از یک موج صوتی کم کیفیت (دارای نویز) یک موج صوتی با کیفیت بالاتر به دست آورد. همچنین این حوزه کاربرد فراوانی برای سرگرمی و استفاده در فضاهای مجازی دارد. برای

مثال افراد می‌توانند صدای خود را به صدای افراد معروف یا صداهای ترسناک یا خنده‌دار تبدیل نمایند. بنابراین وجود بستری برای ارائه سرویس تبدیل صدا بسیار مفید و پرترفدار خواهد بود و به راحتی می‌توان از آن استفاده تجاری نمود.

۱-۳ انواع تبدیل صدا

مدل‌های تبدیل صدا با توجه به روشی که استفاده می‌کنند در نهایت به صورت صفر شات^۱ یا تک شات^۲ انجام می‌شوند. تک شات به معنای آن است که مدل در فرایند یادگیری باید گوینده مقصد را دیده باشد که نوعی محدودیت به حساب می‌آید. ولی مدل‌های صفر شات این محدودیت را ندارند. همچنین برخی مدل‌ها قابلیت تبدیل صدا به صورت بی‌درنگ^۳ را دارند. در این پروژه ما از یک مدل تک شات استفاده کردیم که قابلیت بی‌درنگ بودن را نیز ندارد.

۱-۴ اهداف پژوهش

هدف از این پروژه پیاده سازی بستری برای ارائه سرویس تبدیل صدا می‌باشد که با توجه به کیفیت آن می‌توان در حوزه‌های گوناگون از جمله صنعت دوبله آن را به کار برد و تنها لازم است یک دیتاست از صدای گوینده یا صدای پیشه موردنظر جمع‌آوری نموده و مدل را با آن آموزش دهیم. ما با توجه به امکاناتی که در اختیار بود فرایند یادگیری را بر روی داده‌های صوتی خانم صارمی (گوینده) انجام دادیم.

۱-۵ ساختار پایان‌نامه

این پایان‌نامه در شش فصل به شرح زیر ارائه می‌شود. در فصل دوم به شرح مختصر دو مدل بسیار مهم در حوزه پردازش گفتار می‌پردازیم که از این دو مدل در مقاله‌ی اصلی مورد استفاده در این پایان‌نامه استفاده شده است. در فصل سوم به شرح ساختار و عملکرد مدل FreeVC که مبنای اصلی این پروژه است می‌پردازیم. در فصل چهارم به نحوه‌ی آموزش و استنتاج مدل پرداخته شده است. در فصل پنجم هم پیاده‌سازی سرویس با استفاده از Django و ساختار پایگاه داده را شرح داده‌ایم. فصل ششم به جمع‌بندی کارهای انجام شده در این پژوهش و ارائه‌ی پیشنهادهایی برای انجام کارهای آتی خواهد پرداخت.

¹Zero Shot

²One Shot

³Real Time

فصل ۲

کارهای پیشین

در این پروژه از یک مدل تبدیل صدا به نام FreeVC استفاده شده است. این مدل خود از مدل‌های دیگری درون ساختار خود استفاده کرده است که در این فصل به بررسی مختصر این مدل‌ها می‌پردازیم و در فصل بعد مدل FreeVC را شرح خواهیم داد.

۱-۲ HiFi GAN

یکی از مدل‌هایی که در زمینه تولید گفتار نتایج خیلی خوبی گرفته است مدل HiFi GAN می‌باشد. این مدل که در سال ۲۰۲۰ ارائه شده است هم از لحاظ محاسباتی و هم از لحاظ کیفیت و کارایی نسبت به مدل‌های پیشین خود عملکرد بهتری دارد. HiFi GAN یک مدل تولید کننده^۱ می‌باشد که بر اساس یادگیری تخصصی^۲ کار می‌کند به این صورت که یک تولید کننده (Generator) و دو جدا کننده (Discriminator) در مقابل هم قرار می‌گیرند که در زیر به شرح عملکرد هریک می‌پردازیم:

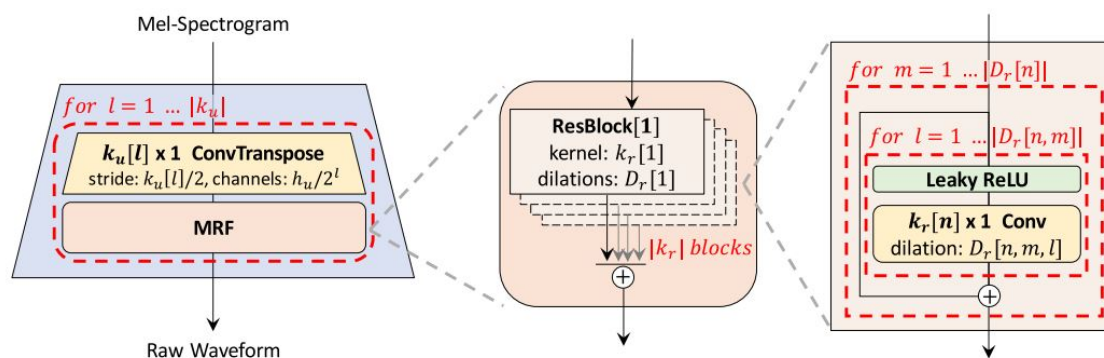
۱-۱-۲ Generator

شبکه تولید کننده در HiFi GAN بر اساس CNN^۳ می‌باشد. ساختار این تولید کننده در شکل ۱-۲ آمده است.

^۱Generative

^۲Adversarial

^۳Convolutional Neural Network



شکل ۲-۱: ساختار Generator در HiFi GAN

ورودی این ماژول، اسپکتروگرام مل^۴ می‌باشد و به عنوان خروجی یک موج صوتی می‌دهد. بر روی هر اسپکتروگرام ورودی تعدادی تبدیل Convolution اجرا می‌شود و پس از هر Convolution یک ماژول MRF^۵ قرار دارد که شامل تعدادی بلاک^۶ می‌باشد که به طور موازی اجرا می‌شوند که هر یک مقداری را به حاصل Convolution اضافه می‌کند.

۲-۱-۲ Discriminator

سیگنال صوتی یک گفتار شامل تعداد زیاد موج سینوسی با طول موج‌های متفاوت است. یکی از مسائلی که در پردازش گفتار با آن روبرو هستیم شناسایی تمام الگوهای متناوب موجود در یک سیگنال صوتی می‌باشد. همچنین برای مدل کردن هر چه بهتر یک سیگنال صوتی لازم است وابستگی‌های بلند مدت آن را شناسایی کنیم. برای مثال یک واج^۷ ممکن است بیش از ۱۰۰ میلی‌ثانیه طول بکشد که باعث ایجاد همبستگی میان بیش از ۲۲۰۰ نمونه^۸ مجاور می‌شود. در این مدل برای حل این مشکلات دو Discriminator قرار داده شده است که به شرح هر یک از آنان می‌پردازیم.

MPD: Multi-Period Discriminator خود شامل تعدادی جداکننده (Sub-Discriminator)

می‌باشد که هر کدام از آن‌ها تعداد p نمونه صوتی را دریافت می‌کند. هر کدام از این جداکننده‌ها ویژگی‌های متفاوتی از یک سیگنال صوتی را استخراج می‌کنند زیرا با قسمت‌های متفاوتی از این سیگنال مواجه شده‌اند. همانطور که در شکل ۲-۲ مشخص است یک سیگنال صوتی یک بعدی با طول T به یک داده‌ساختار

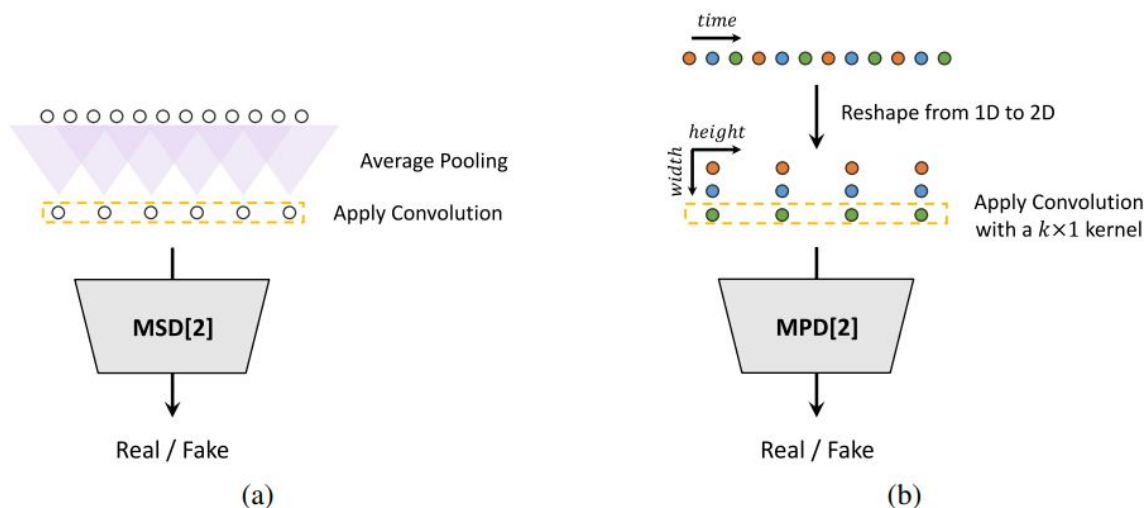
⁴Mel Spectrogram

⁵Multi-Receptive Field Fusion

⁶Residual Block

⁷Phoneme

⁸Sample



شکل ۲-۲: ساختار Discriminator ها در HiFi GAN

دوبعدی با طول T/p و عرض p تبدیل می‌شود. همچنین برای اینکه نمونه‌های متناوب، مستقل از یکدیگر پردازش شوند در هر کانولوشن مقدار Kernel Size را در محور عرض برابر ۱ قرار می‌دهیم.

Multi-Scale Discriminator: همانطور که گفته شد در MPD هر یک از Sub-Discriminator ها نمونه‌های صوتی را به صورت جدا از هم دریافت می‌کند. برای اینکه کل سیگنال صوتی ارزیابی شود یک Discriminator دیگر قرار داده شده‌است. MSD شامل ۳ Sub-Discriminator می‌باشد که هر یک ورودی‌های متفاوتی دریافت می‌کنند:

- Raw Audio
- 2×2 Averaged-pool Audio
- 4×4 Averaged-pool Audio

هرکدام از این Discriminator ها شامل دنباله‌ای از لایه‌های کانولوشنی گروه‌بندی شده و همراه با گام^۹ می‌باشد. لازم به ذکر است که ساختار این Discriminator برگرفته از مدل MelGAN می‌باشد.

^۹Stride

۳-۱-۲ خطای یادگیری

خطای یادگیری تخصصی: در یادگیری تخصصی Generator سعی می‌کند تا سیگنالی تولید کند که Dis-criminator آن را سیگنال صوتی واقعی تشخیص دهد و Discriminator نیز سعی بر آن دارد تا سیگنال واقعی را واقعی و سیگنال تولید شده توسط Generator را غیر واقعی تشخیص دهد. بنابراین خطای یادگیری تخصصی برای تولیدکننده و جداکننده به صورت زیر است:

$$L_{adv}(D; G) = \mathbb{E}_{(x,s)}[(D(x) - 1)^2 + (D(G(s)))^2] \quad (1-2)$$

$$L_{adv}(G; D) = \mathbb{E}_s[(D(G(s)) - 1)^2] \quad (2-2)$$

در عبارات بالا x نشان‌دهنده سیگنال واقعی و s نشان‌دهنده اسپکتروگرام مل ورودی تولید کننده است.

خطای اسپکتروگرام مل: علاوه بر خطای یادگیری می‌توان خطای دیگری اضافه نمود تا بهره‌وری یادگیری و همچنین کیفیت سیگنال خروجی تولیدکننده افزایش یابد. بدین منظور از خطای اسپکتروگرام مل استفاده می‌شود. به گونه‌ای که تولیدکننده سعی می‌کند تا سیگنالی تولید کند که فاصله L_1 میان اسپکتروگرام آن و اسپکتروگرام سیگنال واقعی کمینه شود. این خطا به صورت زیر تعریف می‌شود:

$$L_{Mel}(G) = \mathbb{E}_{(x,s)}[\|\phi(x) - \phi(G(s))\|_1] \quad (3-2)$$

در عبارت بالا ϕ تابع تبدیل کننده موج صوتی به اسپکتروگرام مل آن می‌باشد.

خطای تطابق ویژگی‌ها^{۱۰}: این خطا مربوط به تفاوت میان ویژگی‌های استخراج شده در هر لایه از Discriminator میان سیگنال واقعی و سیگنال تولید شده توسط Generator می‌باشد و به صورت زیر تعریف می‌شود:

$$L_{FM}(G; D) = \mathbb{E}_{(x,s)}\left[\sum_{i=1}^T \frac{1}{N_i} \|D^i(x) - D^i(G(s))\|_1\right] \quad (4-2)$$

در عبارت بالا T نشان‌دهنده تعداد لایه‌های Discriminator و D^i و N_i به ترتیب نشان‌دهنده ویژگی‌ها و تعداد ویژگی‌های لایه i می‌باشد.

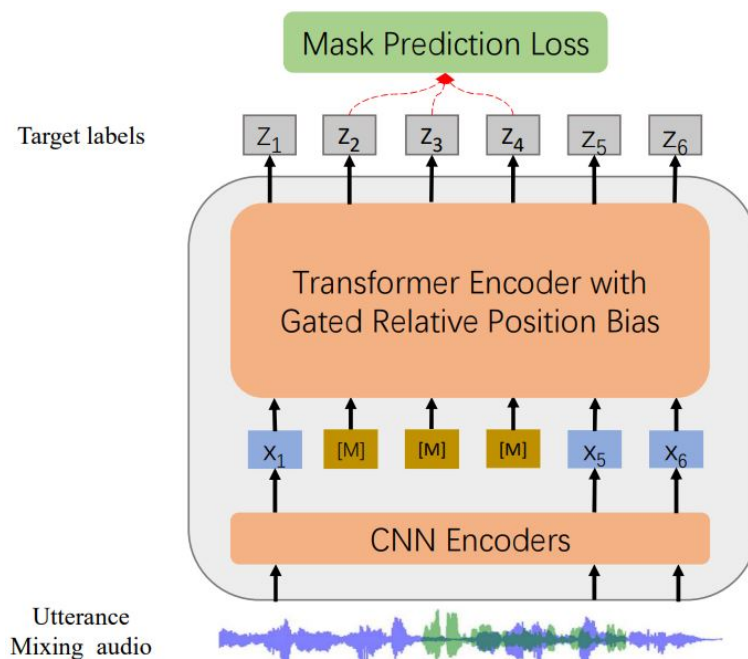
در نهایت خطای کل به شکل زیر درمی‌آید:

$$L_G = L_{Adv}(G; D) + \lambda_{fm} L_{FM}(G; D) + \lambda_{mel} L_{Mel}(G) \quad (5-2)$$

$$L_D = L_{Adv}(D; G) \quad (6-2)$$

که قرار می‌دهیم: $\lambda_{fm} = 2$ و $\lambda_{mel} = 45$

¹⁰Feature Matching



شکل ۲-۳: ساختار WavLM

۲-۲ WavLM

در سال‌های اخیر روش‌های یادگیری خودنظارتی^{۱۱} در حوزه‌های پردازش زبان‌های طبیعی و پردازش گفتار پیشرفت‌های بسزایی رقم زده‌اند. البته در حوزه‌ی پردازش گفتار توجه محققین بیشتر بر روی تشخیص گفتار^{۱۲} معطوف بوده‌است. بنابراین وجود یک مدل از پیش آموزش داده‌شده بر روی حجم زیادی داده برای حوزه‌های مختلف پردازش گفتار بسیار مفید و البته لازم می‌باشد. مدل WavLM که در سال ۲۰۲۲ منتشر شده است، بر روی حجم زیادی از داده بدون برچسب^{۱۳} آموزش دیده است و قابلیت دریافت سیگنال یک گفتار و ارائه نمایش گفتار^{۱۴} را دارد که این نمایش شامل تمام ویژگی‌های مهم گفتار می‌باشد. همچنین WavLM در حذف نویز از سیگنال صوتی هم پیشرفت قابل ملاحظه‌ای رقم زده‌است.

همانطور که در شکل ۲-۳ ملاحظه می‌فرمایید ورودی ماژول WavLM سیگنال صوتی یک گفتار می‌باشد و ابتدا این سیگنال به یک انکودر^{۱۵} CNN داده می‌شود. این انکودر شامل ۷ بلاک کانولوشنی می‌باشد و بعد از هرکدام نیز یک لایه Normalization وجود دارد. خروجی این انکودر به یک انکودر ترانسفورمر داده می‌شود و در نهایت به عنوان خروجی یک نمایش از گفتار، شامل تمام ویژگی‌های مهم اعم

¹¹Self-supervised Learning

¹²Speech Recognition

¹³Unlabeled

¹⁴Speech Representation

¹⁵Encoder

از ویژگی‌های مربوط به گوینده و ویژگی‌های مربوط به محتوای گفتار می‌باشد و البته این خروجی بدون نویز می‌باشد و برای کارهای مختلف پردازش گفتار مانند تشخیص گفتار، تشخیص گوینده و تبدیل صدا قابل استفاده است.

فصل ۳

مدل FreeVC

در این پروژه ما از مدل تبدیل صدای FreeVC^۱ استفاده کرده‌ایم. این مقاله که در سال ۲۰۲۲ منتشر شده است به نتایج بسیار خوبی دست یافته است و برای گفتارهای به زبان فارسی نیز به خوبی عمل می‌کند. در این فصل به شرح ساختار این مدل می‌پردازیم و در فصل بعد گزارش‌های فنی مربوط به آموزش مدل و استنتاج^۲ را ذکر خواهیم کرد.

۳-۱ مدل FreeVC

مدل FreeVC یک مدل تبدیل صدای تک‌شات می‌باشد. رویه‌ی معمولی که در مدل‌های تک‌شات به کار می‌رود اینست که اطلاعات مربوط به محتوا را از گفتار مبدا و اطلاعات مربوط به صدای گوینده را از گفتار مقصد استخراج کرده و با استفاده از یک دکودر^۳ یک موج صوتی جدید بسازیم. بنابراین کیفیت یک مدل به طور کلی به توانایی استخراج اطلاعات مورد نظر از گفتار و توانایی ساخت موج صوتی وابسته است. برخی مدل‌های تبدیل صدا از داده‌های برچسب‌دار استفاده می‌کنند^۴. در این روش معمولاً از یک مدل شناسایی خودکار گفتار^۵ استفاده کرده و اطلاعات مربوط به محتوای گفتار را از این طریق بدست می‌آورند. مشکلی که این روش‌ها دارند برچسب‌گذاری داده‌ها می‌باشد که فرایندی پرهزینه است. بنابراین روش‌های بدون برچسب‌گذاری^۶ طرفدار زیادی پیدا کرده‌اند. از لحاظ ساختار، بسیاری از مدل‌های تبدیل صدا در دو مرحله

^۱FREEVC: TOWARDS HIGH-QUALITY TEXT-FREE ONE-SHOT VOICE CONVERSION

^۲Inference

^۳Decoder

^۴Text-Based Models

^۵Automatic Speech Recognition(ASR)

^۶Text-Free

به انجام می‌رسند. در مرحله اول، اطلاعات آکوستیکی گفتار به اطلاعات مربوط به گوینده مقصد تبدیل می‌شود و در مرحله دوم یک وُکودر^۷ اطلاعات تبدیل شده را به موج صوتی تبدیل می‌کند. مشکلی که در این حالت به وجود می‌آید اینست که این دو مرحله جدا از هم آموزش داده می‌شوند و اطلاعاتی که مدل اول پیش‌بینی می‌کند توزیع احتمال متفاوتی نسبت به داده‌هایی دارد که مدل دوم با آن آموزش داده می‌شود. این مشکل که به آن ناهمخوانی ویژگی‌ها^۸ گفته می‌شود باعث کاهش کیفیت موج صوتی ساخته‌شده می‌شود. مدل FreeVC از داده‌های بدون برچسب استفاده می‌کند و همچنین از طریق اتصال دو مرحله ذکر شده به وسیله بردارهای پنهان^۹ سعی در حل مشکل ناهمخوانی ویژگی‌ها دارد. همچنین علی‌رغم اینکه این مدل تک‌شات می‌باشد ویژگی مثبتی که دارد اینست که لازم نیست مدل با گوینده گفتار مبدا در حین آموزش مواجه شده باشد.

۲-۳ ساختار مدل FreeVC

مدل FreeVC شامل اجزای: انکودر پیشین^{۱۰}، انکودر پسین^{۱۱}، انکودر گوینده^{۱۲}، دکودر و جداکننده می‌باشد که در زیر به شرح عملکرد هریک می‌پردازیم.

۱-۲-۳ انکودر پیشین

انکودر پیشین همانطور که در شکل ۱-۳ مشخص است، خود شامل ۳ بخش می‌باشد: ماژول WavLM، یک گلوگاه اطلاعاتی^{۱۳} و یک ماژول نرمال‌کننده^{۱۴}. ماژول WavLM و گلوگاه، مسئول استخراج اطلاعات مربوط به محتوای گفتار به صورت توزیع $N(z' : \mu_\theta, \sigma_\theta^2)$ می‌باشند.

ماژول WavLM به عنوان ورودی یک موج صوتی دریافت کرده و یک بردار ۱۰۲۴ بعدی از ویژگی‌های SSL^{۱۵}، شامل اطلاعات مربوط به محتوای گفتار و اطلاعات مربوط به گوینده گفتار را به عنوان خروجی می‌دهد (x_{ssl}). به دلیل اینکه نیازی به اطلاعات مربوط به گوینده نداریم باید این اطلاعات را از خروجی WavLM از بین ببریم. به همین دلیل این بردار به یک گلوگاه اطلاعاتی داده می‌شود و بردار ۱۰۲۴ بعدی به

⁷Vocoder

⁸Feature Mismatch

⁹Latent Vectors

¹⁰Prior Encoder

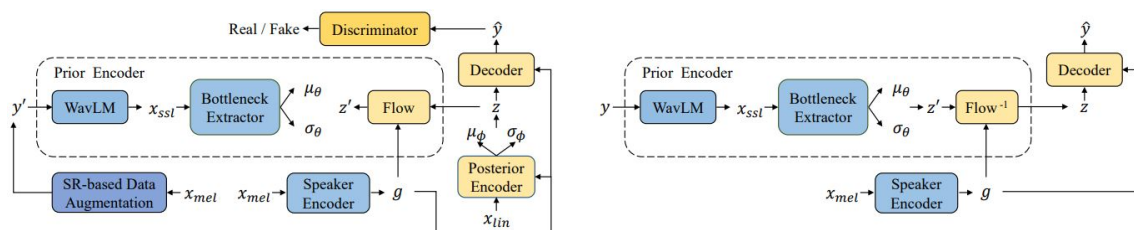
¹¹Posterior Encoder

¹²Speaker Encoder

¹³Information Bottleneck

¹⁴Normalizing Flow

¹⁵Self-Supervised Learning



شکل ۳-۱: سمت راست: استنتاج، سمت چپ: یادگیری

یک بردار با بعد بسیار کوچکتر (d) تبدیل می‌شود.^{۱۶} به دلیل کاهش بعد زیادی که رخ می‌دهد به این قسمت گلوگاه می‌گویند که باعث می‌شود تا اطلاعات مربوط به گوینده و همچنین نویزها و اطلاعات اضافی از بین بروند. در نهایت این بردار d بعدی به دو بردار d بعدی μ_θ و σ_θ تبدیل می‌شود. سپس برای اینکه بر پیچیدگی توزیع احتمال بدست آمده بیافزاییم، این توزیع را به یک ماژول Normalizing Flow می‌دهیم که شامل تعدادی تبدیل آفین می‌باشد. بدین صورت توزیع احتمال ساده‌ی ابتدایی (نرمال) به یک توزیع پیچیده‌تر تبدیل می‌شود.

۳-۲-۲ انکودر گوینده

این انکودر وظیفه استخراج اطلاعات مربوط به گوینده را دارد. در این مقاله از دو نوع انکودر برای این امر استفاده شده: یک انکودر از پیش آموزش داده شده که مخصوص شناسایی گوینده^{۱۷} می‌باشد. یک بار هم از یک انکودر آموزش داده نشده استفاده شده است که معماری آن براساس LSTM^{۱۸} می‌باشد و باید همراه سایر اجزای مدل آموزش داده شود. در این پروژه ما از انکودر دوم استفاده کردیم تا همه اجزا با یکدیگر آموزش داده شوند.

۳-۲-۳ انکودر پسین

این انکودر که فقط در فرایند یادگیری حضور دارد به عنوان ورودی اسپکتروگرام خطی موج صوتی و خروجی انکودر گوینده را دریافت کرده و یک توزیع احتمال نرمال با پارامترهای μ_ϕ و σ_ϕ به عنوان خروجی می‌دهد. خروجی این انکودر به دکودر داده می‌شود تا از آن برای ساخت موج صوتی استفاده کند.

در ادامه همین فصل توزیع خواهیم داد که چگونه حضور انکودر پسین در کنار انکودر پیشین باعث

¹⁷Speaker Verification

¹⁸Long Short-Term Memory

برطرف کردن مشکل ناهمخوانی ویژگی‌ها می‌شود.

۴-۲-۳ دکودر و جداکننده

در این مدل از دکودر به کار رفته در مدل HiFi-GAN استفاده شده است. همانطور که در فصل قبل توضیح دادیم این دکودر شامل چند عملیات کانولوشنی می‌باشد که بعد از هرکدام یک ماژول MRF^{۱۹} قرار دارد.

به عنوان جداکننده^{۲۰} هم از جداکننده MPD^{۲۱} به کار رفته در مدل HiFi-GAN استفاده شده است که توضیحات مربوط به آن را می‌توانید در فصل قبل ملاحظه فرمایید.

۵-۲-۳ داده‌افزایی

یکی از مشکلاتی که بسیاری از مدل‌های یادگیری ماشین با آن مواجه هستند کمبود داده می‌باشد که منجر به یادگیری ناقص می‌شود. همچنین در چنین مدل‌هایی که یک گلوگاه اطلاعاتی دارند، تعیین اندازه‌ی این گلوگاه بسیار مهم است به گونه‌ای که اگر اندازه آن خیلی کم باشد مدل اطلاعات مهمی را ممکن است از دست بدهد و اگر اندازه آن بسیار بزرگ باشد اطلاعات ناخواسته نیز وارد می‌شوند. به جای آنکه اندازه گلوگاه را با دقت تنظیم کنیم، می‌توانیم حجم دادگان خود را با استفاده از روش‌های داده‌افزایی، افزایش داده و مدل در فرایند یادگیری به تدریج بیاموزد که اطلاعات مهم را حفظ و اطلاعات زائد را حذف کند.

در این مدل برای داده‌افزایی از روش تغییر اندازه اسپکتروگرام^{۲۲} استفاده شده است.

در این روش همانطور که در شکل ۲-۳ ملاحظه می‌فرمایید اسپکتروگرام مل یک موج صوتی را دریافت کرده و با یک ضریب r آن را در راستای محور زمان یا محور فرکانس فشرده یا کشیده می‌کنیم. اگر $r > 1$ اسپکتروگرام، کشیده می‌شود و باید بخش اضافی آن را *cut* کرد. و اگر $r < 1$ اسپکتروگرام، فشرده می‌شود و باید جای خالی را *pad* کرد.

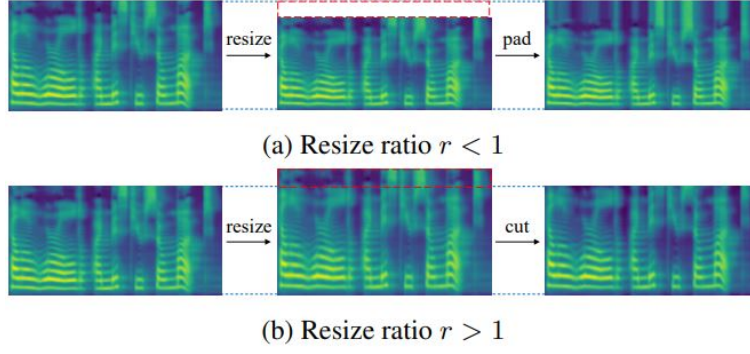
بدین ترتیب می‌توان حجم دادگانی را که در اختیار داریم افزایش داریم و در نهایت کیفیت خروجی مدل را بهبود ببخشیم.

¹⁹Multi-Receptive Field Fusion

²⁰Discriminator

²¹Multi-Period Discriminator

²²Spectrogram Resize



شکل ۳-۲: داده‌افزایی مبتنی بر تغییر اندازه اسپکتروگرام

۳-۳ خطای مدل

خطای مدل FreeVC بسیار به خطای مدل HiFi-GAN شباهت دارد. خطای جداکننده که در دو مدل عیناً مثل هم است:

$$L_{adv}(D) = \mathbb{E}_{(y,z)} [(D(y) - 1)^2 + (D(G(z)))^2] \quad (۱-۳)$$

در مدل HiFi-GAN همانطور که فصل قبل توضیح داده شد برای تولید کننده ۳ خطا داریم:

$$L_{adv}(G; D) = \mathbb{E}_z [(D(G(z)) - 1)^2] \quad (۲-۳)$$

$$L_{Mel}(G) = \mathbb{E}_{(y,z)} [\| \phi(y) - \phi(G(z)) \|_1] \quad (۳-۳)$$

$$L_{FM}(G; D) = \mathbb{E}_{(y,z)} \left[\sum_{i=1}^T \frac{1}{N_i} \| D^i(y) - D^i(G(z)) \|_1 \right] \quad (۴-۳)$$

در این مدل یک خطای دیگر هم داریم:

$$L_{kl} = \log q_\phi(z \mid x_{lin}) - \log p_\theta(z \mid c) \quad (۵-۳)$$

که در عبارت بالا c اطلاعات مربوط به محتوای گفتار می باشد و داریم:

$$q_{\phi}(z | x_{lin}) = N(z; \mu_{\phi}, \sigma_{\phi}^2) \quad (6-3)$$

$$p_{\theta}(z | c) = N(z'; \mu_{\theta}, \sigma_{\theta}^2) | \det \frac{\partial z'}{\partial z} | \quad (7-3)$$

این خطا در واقع دیورژانس کولباک-لیبلر می باشد. این دیورژانس دو توزیع احتمال دریافت کرده و می گوید این دو توزیع از لحاظ اطلاعاتی که در اختیار می گذارند چقدر با یکدیگر تفاوت دارند. در اینجا ما دو توزیع احتمال پیشین و پسین را به این دیورژانس داده ایم در نتیجه با کمینه کردن این خطا، توزیع احتمالی که انکودر پیشین پیش بینی می کند و توزیع احتمالی که دکودر با آن آموزش داده می شود به یکدیگر نزدیک خواهند شد در نتیجه مشکل ناهمخوانی ویژگی ها تا حدودی برطرف می شود.

در نهایت خطای کل مدل به صورت زیر درمی آید:

$$L(D) = L_{Adv}(D) \quad (8-3)$$

$$L(G) = L_{Adv}(G) + L_{Mel}(G) + L_{FM}(G) + L_{kl}(G) \quad (9-3)$$

فصل ۴

گزارش‌های فنی آموزش و استنتاج مدل

در این فصل به شرح چگونگی آموزش مدل و استنتاج آن می‌پردازیم و تنظیمات مربوط به آن را بیان می‌کنیم. نکته‌ی لازم به ذکر آنست که این مدل به ۳ صورت قابل استفاده‌است:

- بدون داده‌افزایی w/o sr
- استفاده از انکودر گوینده‌ی آموزش ندیده و همراه با داده‌افزایی (FreeVC-s)
- استفاده از انکودر گوینده‌ی از پیش آموزش داده و همراه با داده‌افزایی (FreeVC)

نتایج ارزیابی هرکدام از این ۳ حالت را در پیوست آ می‌توانید ملاحظه فرمایید. همچنین لازم به ذکر است که ما در این پروژه از حالت دوم (FreeVC-s) استفاده کردیم.

۴-۱ تنظیمات مدل

پس از دانلود مدل از صفحه گیت‌هاب^۱ می‌توان به فایل‌های آن دسترسی پیدا کرد. در پوشه configs سه فایل آمده است و باتوجه به اینکه ما از دومین حالت مدل استفاده کردیم فایل freevc-s را باز می‌کنیم. حال می‌توان تنظیمات مربوط به مدل را مشاهده و در صورت نیاز تغییر داد. تنها موردی که تغییر داده‌شد مقدار batch size می‌باشد که به دلیل کمبود حافظه مجبور شدیم آن را کاهش دهیم و برابر ۸ قرار دهیم. همچنین همانطور که ملاحظه می‌فرمایید مقدار epochs برابر ۱۰۰۰۰ می‌باشد که آن را تغییر ندادیم ولی

^۱github.com/olawod/freevc

با توجه به checkpoint هایی که در هر مرحله بدست می‌آمد تقریباً پس از گذشت ۲۰۰۰ epoch تغییر چندانی ملاحظه نشد. بنابراین احتمالاً با کاهش این مقدار مدل باز هم به خوبی کار کند. همچنین باید فایل‌های مربوط به Generator و WavLM را از آدرس ذکر شده در صفحه گیت‌هاب دانلود کرده و در پوشه مورد نظر قرار دهیم.

۲-۴ آموزش مدل

ابتدا یک پوشه به نام dataset ایجاد می‌کنیم. در این پوشه نیز به تعداد افرادی که از صدایشان داده داریم پوشه ساخته و نام هر شخص را بر روی یک پوشه می‌گذاریم. برای مثال اگر از ۳ گوینده به نام‌های spk۱ و spk۲ و spk۳ داریم، درون پوشه dataset سه پوشه به نام‌های spk۱، spk۲ و spk۳ می‌سازیم. ما چون از یک گوینده (خانم صارمی) داده‌ی صوتی داشتیم ۱ پوشه به نام saremi ساختم و دادگان صوتی را در این پوشه قرار دادیم. (لازم به ذکر است تمام فایل‌های صوتی باید ترخ نمونه‌برداری^۲ ۱۶KHz داشته باشند). ابتدا قبل از شروع فرایند آموزش باید مدل را آماده کنیم. اگر می‌خواهیم از انکودر گوینده‌ی از پیش آموزش داده شده استفاده کنیم باید تکه کد زیر را اجرا کنیم:

```
# run this if you want to use pretrained speaker encoder
CUDA_VISIBLE_DEVICES=0 python preprocess_spk.py
```

شکل ۴-۱: یادگیری با انکودر از پیش آموزش دیده

اگر نمی‌خواهیم از داده‌افزایی استفاده کنیم باید تکه کد زیر را اجرا کنیم:

```
# run this if you want to train without SR-based augmentation
CUDA_VISIBLE_DEVICES=0 python preprocess_ssl.py
```

شکل ۴-۲: یادگیری بدون داده‌افزایی

و اگر می‌خواهیم از داده‌افزایی استفاده کنیم ابتدا باید در پوشه اصلی یک پوشه به نام training_set بسازیم و سپس تکه کد زیر را اجرا کنیم:

در پوشه filelists سه فایل متنی به نام‌های train و test و val وجود دارد که در هر کدام باید مسیر فایل‌های صوتی مربوط به یادگیری یا ارزیابی با تست را قرار دهیم. بنابراین اگر فایل صوتی با نام file_name.wav از گوینده با نام spk۱ را می‌خواهیم در یادگیری شرکت دهیم باید در فایل train وارد

^۲Sample Rate

```
CUDA_VISIBLE_DEVICES=1 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 68 --max 72
CUDA_VISIBLE_DEVICES=1 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 73 --max 76
CUDA_VISIBLE_DEVICES=2 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 77 --max 80
CUDA_VISIBLE_DEVICES=2 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 81 --max 84
CUDA_VISIBLE_DEVICES=3 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 85 --max 88
CUDA_VISIBLE_DEVICES=3 python preprocess_sr.py --in_dir "dataset/" --wav_dir "training-set/" --ssl_dir "training-set/" --min 89 --max 92
```

شکل ۴-۳: یادگیری با داده‌افزایی

کنیم: training_set/spk\file_name.wav اگر می‌خواستیم این فایل را در ارزیابی (validation) شرکت دهیم همین عبارت را باید در فایل val قرار می‌دادیم. لازم به ذکر است که در هر خط از این فایل باید مسیر یک فایل مشخص شود. با توجه به این که تعداد فایل‌ها زیاد می‌باشد وارد کردن مسیر همه آن‌ها به صورت دستی کاری طاقت‌فرسا است بنا بر این می‌توان از تکه کد زیر استفاده کرد:

```
import os
# folder path
dir_path = r'training_set/'

# list to store files
res = []
# Iterate directory
for path in os.listdir(dir_path):
    # check if current path is a file
    if os.path.isfile(os.path.join(dir_path, path)):
        f = open("filelists/train.txt", "a")
        f.write("training_set/" + path + "\n")
        f.close()
```

شکل ۴-۴: وارد کردن آدرس فایل‌های صوتی

در نهایت فایل متنی باید به شکل زیر در آید:

```
training_set/saremi/559.wav
training_set/saremi/13.wav
training_set/saremi/107.wav
```

شکل ۴-۵: فایل متنی

۳-۴ استنتاج

ابتدا باید آخرین checkpoint حاصل از فرایند یادگیری را گرفته و در پوشه f_0 قرار دهیم و نام این checkpoint را ckpt.pth فرض می‌کنیم. یک فایل صوتی از گوینده مقصد نیاز داریم. نام این فایل را target_speaker.wav در نظر می‌گیریم و این فایل را در پوشه f_1 قرار می‌دهیم.

حال فرض کنید می‌خواهیم فایل صوتی source_speaker.wav را که در پوشه f_2 قرار دارد به صدای گوینده target_speaker تبدیل کنیم. بدین منظور باید در فایل convert.txt که در پوشه اصلی قرار دارد عبارت زیر را وارد کنیم:

```
name|f2/source_speaker.wav|f1/target_speaker.wav
```

شکل ۴-۶: فایل convert.txt

که name هر عبارت دلخواه می‌تواند باشد و نام فایل خروجی برابر آن خواهد بود. همچنین در فایل convert.txt می‌توان چند خط وارد کرد و مدل همه‌ی آن‌ها را به ترتیب اجرا خواهد کرد.

حال برای اجرای فرایند استنتاج اگر از FreeVC-s استفاده می‌کنیم باید تکه کد زیر را اجرا کرد:

```
#!/bin/bash
echo "inference begins!"
CUDA_VISIBLE_DEVICES=0 python3 convert.py --hfile configs/freevc-s.json --ptfile f0/ckpt.pth --txtpath convert.txt --outdir outputs/freevc-s
```

شکل ۴-۷: استنتاج

و اگر از FreeVC استفاده می‌شود باید تکه کد زیر را اجرا نمود:

```
#!/bin/bash
echo "inference begins!"
CUDA_VISIBLE_DEVICES=0 python3 convert.py --hfile configs/freevc.json --ptfile f0/ckpt.pth --txtpath convert.txt --outdir outputs/freevc
```

شکل ۴-۸: استنتاج

در نهایت فایل خروجی در پوشه‌ی output/freevc یا output/freevc-s قرار خواهد گرفت

فصل ۵

پیاده‌سازی سرویس تبدیل صدا

پس از آموزش مدل و آشنایی با نحوه استنتاج لازم است تا بستری فراهم شود که سرویس تبدیل صدا را ارائه دهد. بدین منظور ما از بات تلگرام استفاده کردیم. پس از ساخت بات تلگرامی لازم بود تا یک سیستم Backend ساخته شود و بات از طریق این سیستم با مدل تبدیل صدا ارتباط برقرار کند. برای پیاده‌سازی Backend از Django و REST Framework استفاده نمودیم.

همچنین به منظور استفاده تجاری از این بات، لازم بود تا یک پایگاه داده نیز برای این سرویس در نظر گرفته شود. بدین منظور از MongoDB استفاده کردیم و در ادامه فصل به شرح آن خواهیم پرداخت.

۵-۱ پایگاه داده

برای پایگاه داده‌ی این سرویس دو مجموعه^۱ در نظر گرفته‌ایم. یک مجموعه user که اطلاعات هر کاربر را ذخیره می‌کنیم. در حال حاضر برای این مجموعه دو field در نظر گرفته شده‌است: tg_id که آیدی تلگرام شخص را ذخیره می‌کند و join_date که تاریخ اولین استفاده کاربر از بات را نشان می‌دهد.

مجموعه دومی که در این پایگاه داده داریم user_account می‌باشد که حساب کاربر را نشان می‌دهد. هر بار که کاربر عملیات تبدیل صدا یا خرید سکه را انجام می‌دهد یک سند^۲ به این مجموعه اضافه می‌شود. این مجموعه سه field دارد: tg_id که آیدی تلگرام کاربر را مشخص می‌کند، coins که تعداد سکه‌های کم شده یا اضافه شده به حساب کاربر را نشان می‌دهد (علامت منفی به معنای کم شدن است) و cause که علت این افزایش یا کاهش سکه را نشان می‌دهد و ۳ حالت دارد: initial_charge که

^۱Collection

^۲Document

ارسال می‌شود و بات بار ارسال یک درخواست GET به url داده شده فایل را دریافت می‌کند و برای کاربر ارسال می‌کند. پس از ارسال فایل برای کاربر یک سند برای مجموعه‌ی حساب کاربر در پایگاه داده ارسال می‌شود که این سند حاوی اطلاعات زیر است:

```
_id: ObjectId('6475d106cc8e223499b8637f')
tg_id: "kas..."
coins: -1
cause: "vc_operation"
```

برای خرید سکه کافی است از بات یک درخواست POST به `api/buy/id/coins` ارسال شود که `id` آیدی تلگرام کاربر و `coins` تعداد سکه‌هایی است که می‌خواهد بخرد می‌باشد. کاربر کافی است از منوی خرید سکه یک گزینه را انتخاب کند و خود بات این درخواست را ارسال می‌کند و در جواب لینک درگاه پرداخت داده می‌شود که برای کاربر ارسال می‌شود. در صورت موفقیت آمیز بودن پرداخت یک سند مانند شکل زیر به پایگاه داده ارسال می‌شود و در غیر این صورت کاربر پیغام خطا دریافت می‌کند.

[illegible]

هنگامی هم که کاربر برای اولین بار از بات استفاده می‌کند با کلیک برروی دکمه start بیست سکه به حساب او اضافه می‌شود و پیغام شروع به او ارسال می‌شود:

```
_id: ObjectId('6475c737fc29c943b7da49b5')
tg_id: "kappa-59"
coins: 20
cause: "initial charge"
```

فصل ۶

جمع‌بندی

در این پایان‌نامه به شرح ساختار و عملکرد یک مدل تبدیل صدا مبتنی بر یادگیری عمیق پرداخته شد که بر اساس مدل FreeVC می‌باشد و در نهایت به پیاده‌سازی بستری برای ارائه این سرویس و استفاده تجاری از آن پرداخته شد.

البته این سیستم قابلیت گسترش نیز دارد. به راحتی با جمع‌آوری یک دیتاست^۱ از صدای یک گوینده می‌توان مدل را آموزش داد و گویندگان مختلف را به این سرویس اضافه نمود که قطعا بر جذابیت آن خواهد افزود.

همچنین با توجه به این که برای این سیستم یک Backend تعبیه شده است می‌توان در بسترهای دیگری مانند وب‌سایت یا نرم‌افزار موبایل این سرویس را ارائه داد.

¹Dataset

پیوست آ

مطالب تکمیلی

در این بخش نتایج ارزیابی‌های Subjective و Objective مدل را آورده‌ایم که در زیر می‌توانید ملاحظه فرمایید.

	seen-to-seen		unseen-to-seen		unseen-to-unseen	
	MOS	SMOS	MOS	SMOS	MOS	SMOS
VQMIVC	2.31±0.09	2.10±0.08	1.50±0.08	1.71±0.08	1.49±0.08	1.29±0.05
BNE-PPG-VC	2.80±0.12	2.95±0.12	2.89±0.10	2.83±0.10	3.44±0.08	2.63±0.10
YourTTS	3.46±0.10	3.25±0.09	2.54±0.10	2.50±0.10	2.87±0.09	1.97±0.09
FreeVC	3.99±0.09	3.80±0.09	4.06±0.08	3.77±0.09	4.06±0.08	2.83±0.08
FreeVC (w/o SR)	3.85±0.10	3.50±0.10	3.88±0.08	3.58±0.08	3.97±0.09	2.80±0.09
FreeVC-s	4.01±0.09	3.75±0.09	4.08±0.08	3.68±0.09	4.02±0.09	2.78±0.09
Source	4.32±0.08	-	4.11±0.10	-	4.17±0.09	-

همانطور که در جدول بالا مشاهده می‌کنید بر اساس ارزیابی Subjective که بر اساس نظر مردم می‌باشد در همه‌ی حالات FreeVC عملکرد بهتری داشته است. لازم به ذکر است که MOS نمره‌ای بین ۱ تا ۵ بر اساس کیفیت فایل صوتی خروجی می‌باشد.

SMOS نیز نمره‌ای بین ۱ تا ۵ بر اساس شباهت بین صدای تولید شده و صدای مقصد است. همچنین در حالت سوم که مدل حین یادگیری با صدای مقصد مواجه نشده بود عملکرد خوبی ندارد (از لحاظ شباهت) که نشان‌دهنده‌ی تک شات بودن این مدل است. ولی در حالت دوم که صدای مبدا در حین یادگیری دیده نشده بود همچنان مدل عملکرد خوبی دارد.

در زیر هم ارزیابی Objective مدل آورده شده است که بر اساس نرخ خطای کلمات^۱ و حروف^۲ می‌باشد.

همچنین F₀-PCC که عددی بین ۱- و ۱ است، نشان می‌دهد که مدل FreeVC موفقیت بیشتری

^۱Word Error Rate

^۲Character Error Rate

	WER	CER	F0-PCC
VQMIVC	50.68%	29.61%	0.665
BNE-PPG-VC	6.54%	2.50%	0.718
YourTTS	12.87%	5.70%	0.736
FreeVC	4.35%	1.53%	0.778
FreeVC (w/o SR)	4.92%	1.77%	0.762
FreeVC-s	4.23%	1.46%	0.768

در همخوانی تغییرات فرکانس بین صدای تولید شده و گفتار مبدا دارد.