



واحد تهران جنوب

پایان نامه کارشناسی ارشد (پروپوزال)

فرم شماره ۲

تمامی صفحات طرح تحقیق به صورت تایپ شده تکمیل شود.

عنوان پایان نامه:

فارسی	توسعه یک سیستم تبدیل صوتی کارآمد محاسباتی در تلفن های همراه
انگلیسی	Development of a computationally efficient voice conversion system on mobile phones

مشخصات دانشجو:

نام:	پانیذ	رشته: مهندسی پزشکی	شماره دانشجویی:
نام خانوادگی:	اسلامی تمیجانی	گرایش: بیوالکتریک	۴۰۱۱۴۱۴۰۱۱۱۰۲۶
مجتمع/دانشکده:	دانشکده فنی و مهندسی		
سال تحصیلی اخذ پایان نامه:	۱۴۰۱	ترمه های مشروطی: - تعداد واحدهای گذرانده: ۰ معدل دروس گذرانده شده:	امضاء دانشجو:
نیمسال تحصیلی اخذ پایان نامه:	اول		

کارشناس گروه/مدیر آموزش:

تذکر: اساتید راهنما و مشاور موظف هستند قبل از پذیرش پروپوزال، به سقف ظرفیت راهنمایی و مشاوره خود توجه نموده و در صورت تکمیل نمودن ظرفیت پذیرش، از امضاء این فرم یا در نوبت قرار دادن آن و ایجاد وقفه در کار دانشجویان جدا پرهیز نمایند بدیهی است در صورت عدم رعایت موازین مربوطه، مسئولیت تاخیر در ارائه پروپوزال و عواقب کار، متوجه استاد راهنما خواهد بود.

نام و نام خانوادگی استاد راهنما:	نام و نام خانوادگی استاد مشاور (در صورت لزوم):
امضاء	امضاء

تصویب در شورای گروه تخصصی:	تصویب در شورای پژوهشی مجتمع/دانشکده:
تأیید مدیر گروه	تأیید معاون/مدیر پژوهشی مجتمع/دانشکده
امضاء:	امضاء:
تاریخ:	تاریخ:

# توسعه یک سیستم تبدیل صوتی کارآمد محاسباتی در تلفن های همراه

هدف تبدیل صدا تغییر صدای گوینده منبع به گونه ای است که با حفظ اطلاعات زبانی، صدای گوینده هدف را شبیه صدای گوینده هدف کند. علیرغم پیشرفت سریع الگوریتم های تبدیل صدا در دهه گذشته، بسیاری از آنها هنوز برای عموم پیچیده هستند. با محبوبیت دستگاه های تلفن همراه به ویژه تلفن های هوشمند، برنامه های تبدیل صدای موبایل بسیار مطلوب هستند به طوری که همه می توانند از لذت تقلید صوتی با کیفیت بالا لذت ببرند و افراد مبتلا به اختلالات گفتاری نیز می توانند به طور بالقوه از آن بهره مند شوند. با توجه به محدودیت منابع محاسباتی در تلفن های همراه، نگرانی اصلی این است که بازده زمانی چنین اپلیکیشن موبایلی برای تضمین تجربه کاربر مثبت باشد. در این مقاله، توسعه یک سیستم تبدیل صدای تلفن همراه بر اساس مدل مخلوط گاوسی (GMM) و روش های تاب برداشتن فرکانس وزنی را شرح می دهیم. ما سعی می کنیم با استفاده از بهترین ویژگی های سخت افزاری تلفن های همراه امروزی، مانند محاسبات موازی بر روی چندین هسته و پشتیبانی پیشرفته بردار سازی، کارایی محاسباتی را افزایش دهیم. نتایج ارزیابی تجربی نشان می دهد که سیستم ما می تواند به عملکرد تبدیل صوتی قابل قبولی دست یابد، در حالی که زمان تبدیل برای یک جمله پنج ثانیه ای فقط کمی بیشتر از یک ثانیه در iPhone 7 طول می کشد.

## مقدمه

تبدیل صدا به تغییر صدای یک گوینده که گوینده منبع نامیده می شود، اشاره دارد تا صدا را طوری ایجاد کند که گویی توسط گوینده دیگری به نام گوینده هدف بیان شده است، در حالی که محتوای زبانی بدون تغییر باقی می ماند. بنابراین، ویژگی های صوتی گوینده منبع باید توسط یک سیستم تبدیل صدا برای تقلید از گوینده هدف، بدون تغییر پیام ارسال شده در گفتار، شناسایی و تبدیل شود.

تبدیل صدا پتانسیل زیادی در توسعه کاربردهای صنعتی مختلف دارد، از جمله نه تنها ادغام در سیستم های سنتز متن به گفتار (TTS) برای سفارشی کردن صدا همانطور که می خواهیم [۱]، بلکه نقشی در پزشکی توانبخشی نیز ایفا می کند. مناطق سرگرمی و آموزشی. به عنوان مثال، با تبدیل مصوت های یک گوینده مبتلا به دیس آرتری، یک اختلال حرکتی گفتاری، به فضای مصوت یک گوینده غیر دیس آرتری، درک گفتار به طور قابل توجهی بهبود یافته است [۲]. علاوه بر این، فناوری های تبدیل صدا نیز برای تولید صداهای چند خواننده متنوع با استفاده از یک مدل تبدیل به یک خواننده تک استفاده شده است.

پایگاه داده [۳]، برای ترکیب گفتار احساسی از گفتار خواندن استاندارد (خنثی) [۴] یا تبدیل احساسات منتقل شده در یک گفتار به دیگری [۵]، و تولید نسخه های تصحیح شده عروضی از گفته های زبان آموزان زبان خارجی در آموزش تلفظ به کمک کامپیوتر [۶، ۷]. برعکس، تبدیل صدا ممکن است تهدیدی برای تأیید صوت بلندگو باشد و در نتیجه باعث افزایش ظرفیت ضد جعل سیستم های تأیید بلندگوی معمولی شود [۸]. در مجموع، سزاوار تلاش ما برای مطالعه تبدیل صدا برای اهداف تحقیقات علمی و کاربردی صنعتی است.

به طور کلی، گفتار انسان را می توان به سه جزء تقسیم کرد: محتوای زبان، الگوی طیفی و عروض [۸]. دو عامل آخر تمرکز مطالعات تبدیل صدای فعلی را تشکیل می دهند که در دو سطح واقع شده اند: سطح

فوق‌بخشی، که شامل ویژگی‌های عروضی مانند خطوط فرکانس اساسی، مدت زمان کلمات، واج‌ها، زمان‌بندی، ریتم و سطوح شدت و غیره است. ؛ و سطح سگمنتال شامل شدت گام متوسط، پاسخ فرکانسی دستگاه صوتی، و ویژگی‌های منبع گلوताल. از این نظر، تبدیل صدا با هدف نگاشت ویژگی‌های طیفی و عروضی گوینده مبدا به گوینده هدف به منظور اصلاح فردیت صدا است. چنین وظیفه نگاشت ویژگی معمولاً به عنوان یک مشکل یادگیری نظارت شده در ادبیات فرموله می‌شود که به مقدار معینی از داده‌های گفتاری در دسترس هم از گوینده منبع و هم از گوینده هدف برای آموزش نیاز دارد. تفاوت کلیدی بین انواع مختلف مطالعات در مدل‌های نقشه‌برداری خاص اتخاذ شده، عمدتاً شامل روش‌های مبتنی بر کتاب کد، آماری و مبتنی بر شبکه‌های عصبی مصنوعی است که در یک بررسی اخیر به طور کامل بررسی شده‌اند [۹]. در ادامه، ما فقط برخی از رویکردهای نماینده را به طور مختصر مورد بحث قرار می‌دهیم تا عقلانیت روش شناسی خود را توضیح دهیم.

یک سیستم تبدیل صدای معمولی به دو نوع مدل نیاز دارد، یعنی مدل گفتار و مدل تبدیل. پس از بررسی‌های گسترده، مشخص شد که سه ویژگی گفتار مرتبط با فردیت گوینده عبارتند از: طیف متوسط، شکل‌دهنده‌ها و سطح متوسط گام [۶]. در نتیجه، بیشتر سیستم‌های تبدیل صدا تلاش می‌کنند تا پوشش‌های طیفی کوتاه‌مدت و ویژگی‌های عروضی از جمله مقدار زیر و بم، مدت و شدت را تغییر دهند [۹]. بنابراین، قبل از تبدیل صدا، یک مدل گفتار مناسب باید برای تجزیه و تحلیل سیگنال‌های گفتاری ورودی اتخاذ شود تا ویژگی‌های گفتاری مربوطه برای اصلاحات بعدی استخراج شود. علاوه بر این، یک مدل گفتار خوب باید بتواند سیگنال گفتار را از پارامترهای مدل با کیفیت بالا بازسازی کند، ظرفیتی که ما به آن نیاز داریم تا پس از تبدیل، گفتار را برای گوینده هدف ترکیب کنیم. معمولاً مدل گفتار فریم به فریم ساخته می‌شود در حالی که هر فریم یک بخش کوتاه مدت (معمولاً کمتر از ۲۰ میلی ثانیه) را نشان می‌دهد.

مدل‌های گفتاری رایج عبارتند از مدل [10] STRAIGHT، مدل هم‌پوشانی-افزودن گام-همگام [11] (PSOLA)، و مدل نویز به علاوه هارمونیک [12] (HNM)، و غیره STRAIGHT. متعلق به خانواده مدل فیلتر است که فرض می‌کند گفتار با فیلتر کردن سیگنال تحریک با فیلتر مجرای صوتی مستقل از تحریک تولید می‌شود. برای کاهش تداخل بین تناوب سیگنال و پوشش طیفی، STRAIGHT یک تحلیل طیفی تطبیقی را انجام می‌دهد، که می‌تواند جزئیات سطوح فرکانس زمانی را حفظ کند و در عین حال ساختارهای جزئی ناشی از تناوب سیگنال را تقریباً کاملاً حذف کند. از سوی دیگر، روش‌های مبتنی بر سیگنال مانند PSOLA و HNM گفتار را به عنوان ترکیبی از سیگنال منبع و فیلتر پوشش طیفی مدل‌سازی نمی‌کنند، در نتیجه از فرضیات محدودکننده مانند استقلال بین سیگنال منبع و فیلتر اجتناب می‌کنند. معمولاً منجر به کیفیت بالاتر سنتز گفتار می‌شود. به طور خاص، روش‌های PSOLA می‌توانند سیگنال گفتار را در حوزه فرکانس (FD-PSOLA) یا مستقیماً در حوزه زمان (TD-PSOLA) تغییر دهند. به ویژه، TD-PSOLA شکل موج گفتار را به جریانی از سیگنال‌های تحلیل کوتاه‌مدت تنظیم می‌کند که با نرخ همگام گام تنظیم شده‌اند، که یکی از محبوب‌ترین و ساده‌ترین روش‌ها برای اصلاح عروضی با کیفیت بالا است. با این وجود، در جنبه تبدیل صدا، مدل‌های مبتنی بر تجزیه سینوسی گفتار مطلوب‌تر هستند، زیرا چنین مدل‌هایی با دستکاری پارامترهای مدل، تغییرات طیفی و عروضی انعطاف‌پذیر را امکان‌پذیر می‌کنند. به عنوان یک نماینده معمولی، HNM یکی از پرکاربردترین مدل‌ها برای سنتز و اصلاح

گفتار است HNM [12]. سیگنال گفتار را به یک بخش قطعی به عنوان مجموع سینوئیدها با فرکانس های مربوط به زیر و بم و یک بخش تصادفی که با فیلتر کردن نویز گاوسی سفید به دست می آید، تجزیه می کند. یک مطالعه مقایسه ای در [۱۳] نشان می دهد که HNM عملکرد کلی بهتری نسبت به TD-PSOLA و در نتیجه یک مدل گفتاری مبتنی بر ارائه می دهد.

در چنین تجزیه تصادفی هارمونیک به علاوه در سیستم ما به تصویب رسید. برگرفته از مرحله مدل سازی و تحلیل گفتار فوق، ویژگی های خروجی ممکن است مستقیماً مورد استفاده قرار گیرند یا برای نقشه برداری ویژگی بعدی پردازش شوند [۹]. همانطور که قبلاً ذکر شد، در تبدیل صدای فعلی، اکثر روش ها پردازش فریم به فریم را فرض می کنند و بنابراین عمدتاً به ویژگی های گفتار محلی موجود در بخش های کوتاه مدت بستگی دارند. ویژگی های طیفی رایج عبارتند از ضرایب Mel-cepstral (MCCs)، ضرایب پیش بینی خطی (LPCCs) Cepstral، فرکانس طیف خط (LSF) و فرکانس formant، و پهنای باند [۸، ۹]. با توجه به ویژگی های عروزی، ما معمولاً فقط الگوهای  $f_0$  و مدت زمان را در نظر می گیریم. سیستم های تبدیل صدای معمولی فقط تغییرات ساده ای از ویژگی های عروزی را انجام می دهند، به عنوان مثال، عادی سازی میانگین جهانی و واریانس مقادیر  $f_0$  با مقیاس ورود به سیستم [۸، ۱۴]، زیرا عروض یک ویژگی فوق بخشی است که بر اساس بخش بندی منتقل نمی شود، اما از طریق واحدهای بزرگتر، که چالش برانگیزتر است [۱۵]. در نتیجه، بیشتر تحقیقات بر روی نقشه برداری برای ویژگی های طیفی تمرکز دارند. از نظر ریاضی، تبدیل صدا عبارت است از یادگیری تابع نگاشت  $f(\cdot)$  از ویژگی گفتار مبدأ  $X$  به ویژگی گفتار هدف  $Y$  با استفاده از پیکره آموزشی و سپس اعمال این تابع به یک نمونه گفتار منبع نادیده جدید برای تبدیل در زمان اجرا

شروع از کار اصلی مبتنی بر کمی سازی برداری (VQ) و کتاب های کد نگاشت توسعه یافته توسط آبه و همکاران. در سال ۱۹۸۸ [۱۶]، روش های بسیاری برای نگاشت ویژگی های طیفی در ادبیات ارائه شده است. اگرچه روش های مبتنی بر نگاشت کتاب کد ساده و از نظر محاسباتی کارآمد هستند، اما به دلیل تولید توالی ویژگی های ناپیوسته، عملکرد ضعیفی داشتند [۹، ۱۵]. در حال حاضر، روش های رایج تر عمدتاً در دسته های زیر طبقه بندی می شوند: (۱) ترکیبی از نگاشت های خطی، مانند مدل های مخلوط گاوسی [13] (GMM)، ۱۷، ۱۸ [و مدل های پنهان مارکوف] 19 (HMM)، ۲۰ (2). [یک مدل نگاشت غیرخطی منفرد، شامل رگرسیون ماشین بردار پشتیبان [۲۱]، شبکه های عصبی مصنوعی [22] (ANN)، ۲۳]، و رویکردهای جدیدتر یادگیری عمیق [۲۴، ۲۵]. (۳) رویکردهای نقشه برداری مبتنی بر تاب فرکانس مانند تاب برداشتن فرکانس وزنی [۲۶]، تاب برداشتن فرکانس دینامیکی [۲۷] و تاب برداشتن فرکانس دوخطی [۲۸]. و (۴) روش های نگاشت ناپارامتریک مانند روش مبتنی بر مثال [۲۹، ۳۰]، رگرسیون فرآیند گاوسی [۳۱] و رویکرد هیستوگرام [32] K در روش های مبتنی بر GMM، توزیع بردارهای ویژگی طیفی منبع (هدف) با گروهی از توزیع های نرمال چند متغیره مدل سازی می شوند، و تابع تبدیل معمولاً یک فرم خطی را اتخاذ می کند و بنابراین می تواند مستقیماً با حداقل مربعات حل شود [۱۳، ۱۵]. روش دیگر برای انجام نقشه برداری در GMM این است که بردارهای ویژگی منبع و هدف را با هم ترکیب کنیم تا یک GMM برای بردار تقویت شده  $z = [x^T, y^T]^T$  که GMM مشترک نامیده می شود و بردار هدف نگاشت شده در طول تبدیل می تواند بسازد. با استفاده از بردارهای میانگین و ماتریس های کوواریانس که از آموزش به دست آورده ایم [۳۳]

به دست می آید. در رویکردهای مبتنی بر ANN، شبکه‌های کم عمق و عمیق برای نقشه‌برداری غیرخطی ویژگی‌ها از منبع به مقصد مورد بررسی قرار گرفته‌اند که ممکن است به دلیل ساختار بسیار منعطف و توانایی برازش غیرخطی برجسته به ویژه زمانی که از یک شبکه عصبی عمیق (DNN) استفاده می‌شود [۲۳، ۳۴]. با این وجود، طراحی و آموزش شبکه به تخصص زیادی نیاز دارد، زیرا فرآیند یادگیری به راحتی می‌تواند در حداقل‌های محلی گیر کند و DNN به طور کلی حجم زیادی از داده‌های آموزشی را درخواست می‌کند [۳۵]. مشخص است که روش‌های آماری مانند GMM تمایل به تولید گفتار بیش از حد صاف با کیفیت پایین دارند، اگرچه می‌توانند هویت گوینده را به خوبی تبدیل کنند [۲۷، ۳۰]. در مقابل، تاب فرکانس با هدف نگاشت محور فرکانس طیف بلندگوی منبع به بلندگوی هدف با تاب برداشتن پوشش طیفی منبع، در نتیجه حفظ جزئیات طیفی بیشتر و تولید کلام با کیفیت بالاتر به طور کلی [۳۰، ۳۶] است. با وجود امتیاز با کیفیت بالا، شباهت بین صداها تبدیل شده و هدفمند، یعنی عملکرد تبدیل هویت، در تاب خوردگی فرکانس رضایت بخش نیست، زیرا شکل طیفی به طور کامل اصلاح نشده است [۲۶]. یک روش جایگزین برای تولید گفتار با کیفیت بالا، تبدیل صدای غیرپارامتری مبتنی بر نمونه است. یک نمونه به عنوان بخشی از طیف‌نگار گفتار که فریم‌های متعدد را در بر می‌گیرد، تعریف می‌شود، در حالی که مجموعه وزن‌های ترکیبی خطی یک بردار فعال‌سازی را تشکیل می‌دهند. برای جلوگیری از اثر هموارسازی بیش از حد، بردار فعال‌سازی باید پراکنده باشد، که می‌توان آن را با فاکتورسازی ماتریس غیرمنفی (NMF) با محدودیت پراکنده یا دکانولوشن ماتریس غیرمنفی تخمین زد [۲۹]. این چارچوب بیشتر گسترش یافته است تا شامل ضریب فشرده‌سازی طیفی و یک تکنیک جبران باقیمانده [۳۰] شود. نتایج آزمایش بر روی پایگاه داده VOICES برتری روش مبتنی بر نمونه را در شباهت و کیفیت گفتار تبدیل شده نشان می‌دهد. با این حال، این رویکرد دارای معایب هزینه محاسباتی بالا است که آن را در حال حاضر برای برنامه‌های کاربردی تلفن همراه نامناسب می‌کند. مسلماً GMM هنوز یکی از موفق‌ترین و پرکاربردترین مدل‌ها در سیستم‌های تبدیل صدای عملی است [۸، ۹، ۱۵، ۳۷]. در سیستم تبدیل صدای خود که برای تلفن‌های هوشمند طراحی شده است، ما از روش تاب برداشتن فرکانس وزنی استفاده کردیم که GMM را با تاب برداشتن فرکانس ترکیب می‌کند تا تعادل بهتری بین کیفیت گفتار و شباهت حاصل شود. [26]

در سال‌های اخیر، دستگاه‌های تلفن همراه، به ویژه تلفن‌های هوشمند و تبلت‌ها، در زندگی روزمره ما غالب شده‌اند. یک گزارش نشان می‌دهد که زمان صرف شده روی دستگاه‌های تلفن همراه بین سال‌های ۲۰۱۴ تا ۲۰۱۵ ۱۱۷ □ رشد داشته است، در حالی که استفاده کلی از برنامه‌های تلفن همراه (برنامه) به طور متوسط ۵۸ □ سال به سال افزایش یافته است [۳۸]. با این حال، اگرچه الگوریتم‌های زیادی با مجموعه داده‌های خاص در ادبیات پیشنهاد و تأیید شده است، همانطور که در بالا ذکر شد، تا آنجا که دانش ما وجود دارد، هنوز هیچ سیستم تبدیل صوتی با استفاده از چنین روش‌های پیشرفته‌ای وجود ندارد که بتواند روی دستگاه‌های تلفن همراه، به ویژه در تلفن‌های همراه کار کند. علاوه بر این، از طریق جستجو در فروشگاه اپل و گوگل پلی برای برنامه‌های فعلی با برچسب تغییر/تبدیل صدا، متوجه شدیم که اکثر برنامه‌های تغییر صدای موجود صرفاً سعی می‌کنند برخی از ویژگی‌های طیفی یا عروضی را بدون هیچ هدف خاصی تغییر دهند و فقط به تغییر صدا کمک کنند. صداها ربات مانند را برای سرگرمی تولید کنید، مانند VoiceLab و

Voice Changer. تنها چند برنامه عمداً برای تقلید صدای شخص از پیش تعریف شده دیگر توسعه یافته اند، مانند Trump Voice Changer، که می تواند متن ورودی را به صورت ترامپ مانند بیان کند.

صدا، و برنامه Celebrity Voice Changer Lite، که لیست ثابتی از افراد مشهور را به عنوان اهداف تبدیل صوتی شما ارائه می دهد. با این حال، طبق بررسی های کاربران، این برنامه های کاربردی تلفن همراه (برنامه ها) محدودیت های آشکاری از خود نشان می دهند، به عنوان مثال، گفته های تولید شده معمولاً غیرطبیعی، غیرمشابه و فاقد قابلیت درک به نظر می رسد. علاوه بر این، برخی از برنامه ها برای دسترسی به سرورهای آنلاین نیاز به اتصال اینترنت دارند. مهمتر از آن، هیچ برنامه ای از سفارشی سازی بلندگوهای هدف پشتیبانی نمی کند، بلکه فقط اهداف ثابتی را ارائه می کند، و ویژگی های کلیدی یک سیستم تبدیل صوتی عمومی را از دست می دهد. بنابراین، این به ما انگیزه می دهد تا با اجرای کارآمد یک الگوریتم مبتنی بر GMM، که شامل بهترین استفاده از دستورالعمل های برداری سخت افزار، محاسبات موازی بر روی تلفن های همراه، یک سیستم تبدیل صوتی کامل، آفلاین و بلادرنگ در تلفن های همراه ایجاد کنیم. هسته های متعدد و طراحی خوب معماری نرم افزار. یک برنامه iOS به نام Voichap برای نشان دادن امکان پذیری تبدیل صدا در زمان واقعی با امکانات کامل در دستگاه های آیفون توسعه یافته است.

این مقاله به شرح زیر سازماندهی شده است. در بخش دوم، به طور خلاصه چارچوب الگوریتم زیربنایی برای تبدیل صدا در سیستم خود را معرفی می کنیم. در بخش سوم، ما اصول و تکنیک های اصلی را برای اجرای کارآمد الگوریتم های اصلی در تلفن های همراه شرح می دهیم. بخش چهارم بعدی معماری نرم افزار و توسعه یک برنامه کاربردی کاربردی موبایل برای پلتفرم iOS را ارائه می دهد. بخش V برنامه iOS را نشان می دهد و سپس اثربخشی و کارایی این برنامه به صورت تجربی ارزیابی می شود. در بخش ششم، بحث مختصری در مورد رویکردهای مبتنی بر یادگیری عمیق ارائه می کنیم. در نهایت، مقاله را در بخش هفتم به پایان می رسانیم.

## چارچوب تبدیل صدا

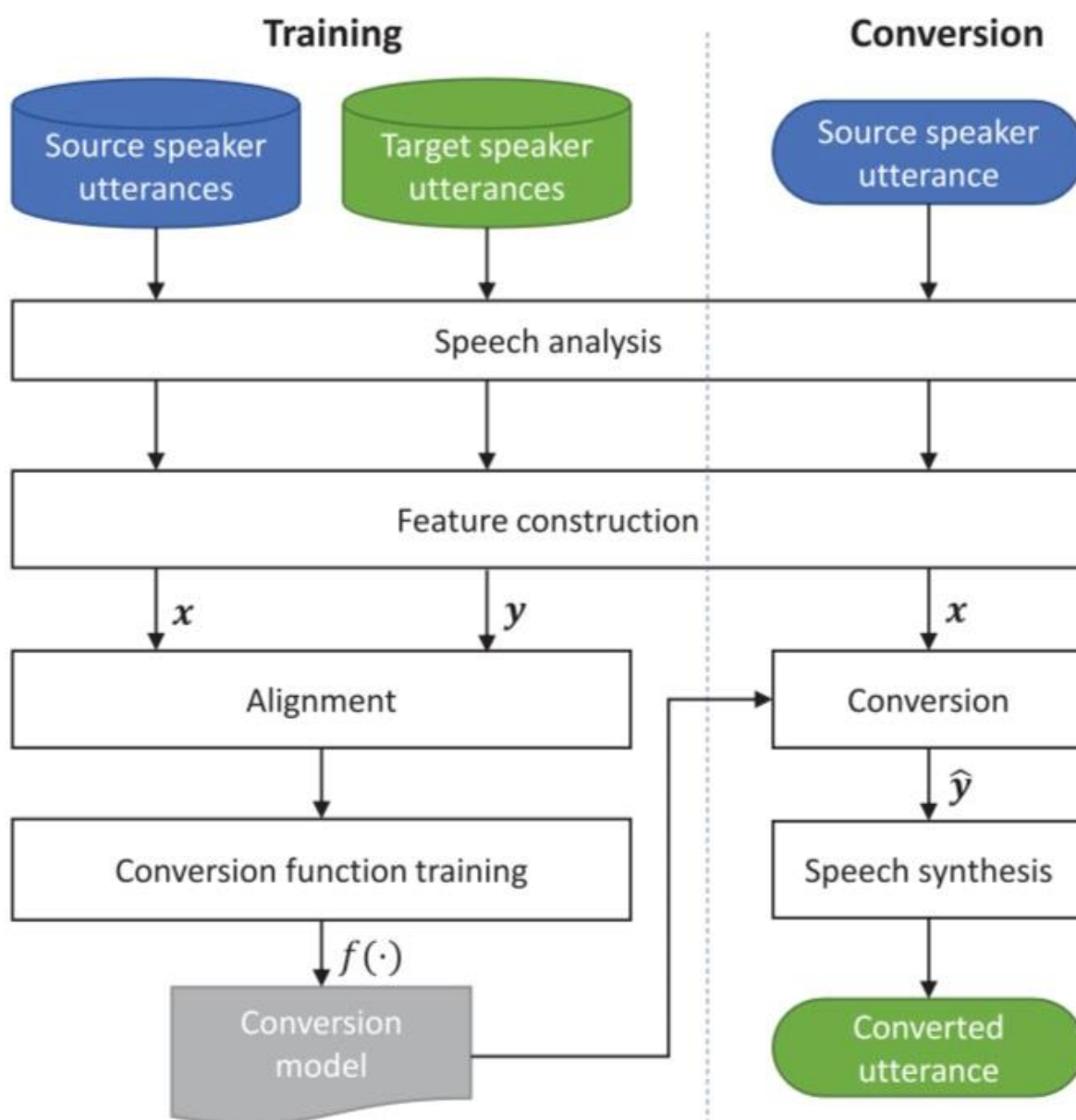
صدا، و برنامه Celebrity Voice Changer Lite، که لیست ثابتی از افراد مشهور را به عنوان اهداف تبدیل صوتی شما ارائه می دهد. با این حال، طبق بررسی های کاربران، این برنامه های کاربردی تلفن همراه (برنامه ها) محدودیت های آشکاری از خود نشان می دهند، به عنوان مثال، گفته های تولید شده معمولاً غیرطبیعی، غیرمشابه و فاقد قابلیت درک به نظر می رسد. علاوه بر این، برخی از برنامه ها برای دسترسی به سرورهای آنلاین نیاز به اتصال اینترنت دارند. مهمتر از آن، هیچ برنامه ای از سفارشی سازی بلندگوهای هدف پشتیبانی نمی کند، بلکه فقط اهداف ثابتی را ارائه می کند، و ویژگی های کلیدی یک سیستم تبدیل صوتی عمومی را از دست می دهد. بنابراین، این به ما انگیزه می دهد تا با اجرای کارآمد یک الگوریتم مبتنی بر GMM، که شامل بهترین استفاده از دستورالعمل های برداری سخت افزار، محاسبات موازی بر روی تلفن های

همراه، یک سیستم تبدیل صوتی کامل، آفلاین و بلادرنگ در تلفن‌های همراه ایجاد کنیم. هسته‌های متعدد و طراحی خوب معماری نرم افزار. یک برنامه iOS به نام Voichap برای نشان دادن امکان‌پذیری تبدیل صدا در زمان واقعی با امکانات کامل در دستگاه‌های آیفون توسعه یافته است.

این مقاله به شرح زیر سازماندهی شده است. در بخش دوم، به طور خلاصه چارچوب الگوریتم زیربنایی برای تبدیل صدا در سیستم خود را معرفی می‌کنیم. در بخش سوم، ما اصول و تکنیک‌های اصلی را برای اجرای کارآمد الگوریتم‌های اصلی در تلفن‌های همراه شرح می‌دهیم. بخش چهارم بعدی معماری نرم افزار و توسعه یک برنامه کاربردی کاربردی موبایل برای پلتفرم iOS را ارائه می‌دهد. بخش V برنامه iOS را نشان می‌دهد و سپس اثربخشی و کارایی این برنامه به صورت تجربی ارزیابی می‌شود. در بخش ششم، بحث مختصری در مورد رویکردهای مبتنی بر یادگیری عمیق ارائه می‌کنیم. در نهایت، مقاله را در بخش هفتم به پایان می‌رسانیم.

## ۱. نمای کلی سیستم تبدیل صدا

یک نمای کلی از یک سیستم تبدیل صدای معمولی، شامل فاز آموزشی و فاز تبدیل در شکل ۱ نشان داده شده است. به طور کلی، در مرحله آموزش، تابع تبدیل  $f(\cdot)$  برای ترسیم بردار ویژگی منبع  $X$  به بردار ویژگی هدف  $Y$  سپس، در طول مرحله تبدیل، بردار ویژگی  $X$  یک گفته منبع جدید است



تصویر ۱: معماری کلی سیستم تبدیل صدا فاتیپیک

ساخته شده و متعاقباً توسط  $\hat{y}=f(x)$  تغییر شکل داده است، که در نهایت برای سنتز گفتار تبدیل شده در صدای گوینده مورد نظر استفاده شد. در اصطلاحات یادگیری ماشین، چنین تبدیل صدا یک مشکل رگرسیونی است، اما به طور کلی یک فرآیند چالش برانگیز است زیرا کیفیت گفتار و تشابه تبدیل تبدیل به هدف معمولاً دو هدف متناقض هستند و باید تعادل مناسبی حاصل شود [۱۵]. در بخش‌های فرعی زیر، مربوط به بلوک‌های شکل ۱، مدل‌سازی گفتار را برای تجزیه و تحلیل گفتار، ساخت ویژگی و یادگیری تابع تبدیل به طور جداگانه شرح می‌دهیم.

## ۲. مدل قطعی و تصادفی



در یک سیستم تبدیل صدا، یک مدل گفتار ابتدا باید برای اهداف سنتز مناسب باشد، به عنوان مثال، وقتی سیگنال گفتاری از پارامترهای مدل بازسازی می‌شود، باید وفاداری را حفظ کند تا از نظر ادراکی از مدل اصلی قابل تشخیص نباشد TD-PSOLA. فوق‌الذکر یک روش قدرتمند برای سنتز گفتار است [۱۱]. با این حال، برای تبدیل صدا مناسب نیست زیرا هیچ مدلی برای سیگنال گفتار در نظر نمی‌گیرد و در نتیجه هیچ دستکاری طیفی را نمی‌توان به راحتی اعمال کرد. بنابراین، سیستم ما از مدل قطعی به علاوه تصادفی بر اساس تجزیه سینوسی گفتار استفاده می‌کند [۳۹]. مشابه HNM کلاسیک [۱۲]، این مدل سیگنال گفتار  $s$  را به عنوان مجموع مجموعه‌ای از سینوسی‌ها با پارامترهای متغیر با زمان و یک جزء نویز مانند نشان می‌دهد که باقیمانده را نشان می‌دهد.

(1)

$$s[n] = \sum_{j=1}^J A_j[n] \cos(\theta_j[n]) + e[n]$$

جایی که قسمت قطعی فقط در قطعات صدا دار ظاهر می‌شود و قسمت تصادفی شامل اجزای سیگنال غیر سینوسی باقی می‌ماند مانند اصطحکاک و صدای تنفس اگرچه پارامترهای مدل در سطح جهانی متغیر هستند، اما می‌توان آن‌ها را در فواصل زمانی کوتاه ثابت در نظر گرفت، و تحلیل محلی سیگنال توسط فریم‌ها را منطقی و آسان‌تر می‌کند، جایی که هر فریم با تعداد ثابتی از نمونه‌های گفتاری مطابقت دارد، مثلاً  $N$ ، در یک فاصله زمانی ۱۰ میلی ثانیه در مورد ما.

پس از تشخیص دامنه‌ها و فازها [۴۰] برای هر نقطه اندازه‌گیری  $k$  مربوط به لحظه زمانی  $kN$ ،  $k \geq 1$ ، شکل موج قطعی در هر لحظه زمانی توسط

(2)

$$A_j[kN + m] = A_j^{(k)} + \frac{A_j^{(k+1)} - A_j^{(k)}}{N} m$$

از  $m = 0, 1, \dots, N-1$  تا جایی که  $A_j^{(k)}$  دامنه هارمونیک  $j$  را در نقطه  $k$  نشان می‌دهد. فازها و فرکانس‌ها توسط پلی سه مرتبه درون یابی می‌شود به شرح زیر:

(3)

$$\theta_j[kN + m] = am^3 + bm^2 + cm + d$$

که در آن پارامترهای  $a$ ،  $b$ ،  $c$ ،  $d$  به صورت بهینه انتخاب می‌شوند روش [۴۱]. اکنون با سینوسی‌های موجود، می‌توانیم بخش قطعی کامل  $d[n]$  را مطابق با آن بدست آورید معادله (۱) و مولفه تصادفی  $e[n]$  می‌تواند باشد

متعاقباً به صورت زیر جدا شد

(4)

$$d[n] = \sum_j^{J(K)} A_j[n] \cos(\theta_j[n])$$

$$e[n] = s[n] - d[n]$$

که در آن  $J(k)$  تعداد هارمونیک ها در قاب  $k$  ام است. کار باقی مانده آنالیز شکل طیفی قدر باقیمانده با روش کدگذاری اعتباری خطی (LPC) است. پس از به دست آوردن مدل قطعی به علاوه تصادفی (۱)، عروض شامل مدت و گام را می توان به طور مستقیم تغییر داد، که در [۴۱] به تفصیل آمده است. برای بازسازی سیگنال از پارامترهای اندازه گیری شده، می توان از تکنیک همپوشانی-افزودن (OLA) برای بازسازی بخش قطعی استفاده کرد، جایی که یک قاب از نمونه های  $2N$  در آن ساخته شده است. هر نقطه اندازه گیری  $k$  که توسط

(5)

$$d[kN + m] = \sum_{j=1}^{J(k)} \left( A_j^{(K)} \cos(w_j^{(k)} + \varphi_j^{(k)}) \frac{N-m}{N} + A_j^{(K+1)} \left( w_j^{(k+1)}(m-N) + \varphi_j^{(K+1)} \right) \frac{m}{N} \right)$$

در مرحله بعد، فریم های طول  $N$  نویز گاوسی سفید در حوزه فرکانس توسط فیلترهای LPC محاسبه شده قبلی برای تولید مولفه تصادفی شکل می گیرند. به این معنی، سیگنال گفتار را می توان با موفقیت از پارامترهای مدل سنتز کرد. نتایج تجربی نشان می دهد که خروجی سیستم از نظر ادراکی از گفتار اصلی قابل تشخیص نیست.

### ۳. فرکانس های طیفی خط دارای ویژگی ساختاری

اگرچه ما یک مدل هارمونیک به علاوه تصادفی ساخته ایم، باید توجه داشت که تبدیل صداها مستقیماً از این پارامترهای مدل بسیار دشوار است، زیرا دامنه ها و فازها پارامترسازی مناسبی از پوشش طیفی هارمونیک برای هدف صدا ارائه نمی کنند. تبدیل. این عمده تاً به دلیل این واقعیت است که تعداد هارمونیک ها متغیر و به طور کلی زیاد است که تبدیل را بسیار پیچیده می کند. بنابراین، همانطور که در شکل ۱ نشان داده شده است، یک مرحله ساخت ویژگی بیشتر مورد نیاز است تا امکان نمایش بهتر گفتار برای هدف تبدیل فراهم شود. در این مطالعه، ضرایب فرکانس های طیفی خطی (LSF) که دارای ویژگی های کوانتیزه سازی و

درون بایی بهتری هستند [۴۲]، به عنوان ویژگی هایی برای تخمین پارامترهای  $\{a_i\}$ ،  $\mu_i$ ،  $i$  مدل مخلوط گاوسی استفاده می شوند. که در زیربخش بعدی معرفی می شود.

با توجه به نمایش تمام قطبی مرتبه  $p$ th از طیف،  $1/A(z)$  ضرایب LSF ریشه های چند جمله ای های مرتبه  $(p + 1)$  هستند که توسط

(6)

$$\begin{aligned} P(z) &= A(z) + z^{-(p+1)} A(z^{-1}) \\ Q(z) &= A(z) - z^{-(p+1)} A(z^{-1}) \end{aligned}$$

که در آن  $P$  یک چند جمله ای پالندرومیک و  $Q$  یک چند جمله ای ضد پالندروم است. توجه داشته باشید که ریشه های  $P(z)$  و  $Q(z)$  به صورت جفت متقارن روی دایره واحد قرار دارند، به این معنی که می توان آنها را کاملاً با فرکانس های مربوط به مکان ریشه ها مشخص کرد و فقط فرکانس های  $p/2$  باید مشخص شوند. برای هر چند جمله ای ذخیره می شود. بنابراین، ویژگی  $LSF$  در مجموع دارای صفات  $p$  است. با توجه به اجرای عملی، تکنیک تکراری مدلسازی تمام قطبی (DAP) استفاده می شود زیرا منجر به اعوجاج کمتری می شود و بنابراین کیفیت ادراکی بهتری را نسبت به پیشینیان خود ارائه می دهد [43]. در تعیین ترتیب بهینه برای فیلترهای تمام قطبی هارمونیک، باید یک معاوضه انجام شود زیرا فیلترهای مرتبه بالا وضوح بالاتر و کیفیت بالاتری را ارائه می دهند، در حالی که فیلترهای مرتبه پایین را می توان با اطمینان بیشتری تبدیل کرد. از طریق آزمون و خطا، فیلترهای تمام قطبی مرتبه چهاردهم  $p = 14$  برای ارائه بهترین نتایج یافت می شوند.

#### ۴. تبدیل صد از طریق تاب برداشتن فرکانس وزنی

از مجموعه آموزشی داده شده، به عنوان مثال، یک پیکره موازی شامل سخنان گوینده منبع و سخنان هدف، بردار ویژگی نهایی  $LSF$  را می توان به ترتیب به صورت  $X$  و  $Y$  بدست آورد. پس از تراز زمانی، استفاده از  $GMM$  برای نشان دادن توزیع ویژگی معمول است. ما می توانیم از یک بردار تقویت شده  $z = [x^T, y^T]^T$  برای ساختن یک  $GMM$  مشترک با مولفه های  $m$  استفاده کنیم.

(7)

$$p(z) = \sum_{i=1}^m \alpha_i N(z; \mu_i, \Sigma_i)$$

جایی که  $m$  و  $E$  به ترتیب بردار و میانگین هستند ماتریس کوواریانس برای مولفه  $i$ th گاوسی و  $(a;)$  مثبت هستند و تا ۱ جمع می شوند  $\mu_i$  و  $\Sigma_i$  را می توان به شکل بلوکی که مربوط به  $X$  و  $Y$  داده شده توسط

(8)

$$\mu_i = \begin{bmatrix} \mu_i^x \\ \mu_i^y \end{bmatrix}, \Sigma_i = \begin{bmatrix} \Sigma_i^{xx} & \Sigma_i^{xy} \\ \Sigma_i^{yx} & \Sigma_i^{yy} \end{bmatrix}$$

که می توان از روی داده ها با استفاده از تکنیک انتظار- حداکثرسازی (EM) تخمین زد. در طول تبدیل، برای بردار ورودی منبع  $x_t$  از قالب  $t$ ، توزیع شرطی  $y_t$  داده شده  $x_t$  دوباره با توزیع گاوسی مخلوط مدل سازی می شود و می توانیم بردار میانگین را به عنوان بردار خروجی هدف پیش بینی شده  $\hat{y}_t$  استفاده کنیم. بنابراین، تابع تبدیل کلاسیک [15]  $GMM$  به صورت زیر فرموله شده است:

(9)

$$F(x_t) = \sum_{i=1}^m \omega_i(x_t) [\mu_i^y + \varepsilon_i^{yx} (\varepsilon_i^{xx})^{-1} (x_t - \mu_i^x)]$$

که در آن  $\omega_i(x_t)$  احتمال خلفی است که LSF داده شده است بردار متعلق به ITH گاوسی است که توسط

(10)

$$\omega_i(x_t) = \frac{\alpha_i N(z; \mu_i, \varepsilon_i)}{\alpha_j N(z; \mu_j, \varepsilon_j)}$$

اگرچه روش سنتی GMM که در بالا توضیح داده شد می‌تواند به شباهت خوبی برای تبدیل صدا دست یابد، کیفیت گفتار تبدیل شده هنوز رضایت‌بخش نیست، عمدتاً به دلیل اثر هموارسازی بیش از حد [۱۵]. در مقابل، روش‌های مبتنی بر تاب‌خوردگی فرکانس می‌توانند از افت قابل توجه کیفیت گفتار جلوگیری کنند، زیرا درجه اصلاح محدود است [۳۵]. بنابراین، ما روشی را که در ابتدا در [۳۶] پیشنهاد شده بود، اتخاذ کردیم، که هدف آن به دست آوردن تشابه تبدیل شده به هدف بالا در عین حفظ کیفیت گفتار، با ترکیب GMM و تاب برداشتن فرکانس است. الهام اصلی این است که میانگین بردارهای LSF هر کلاس آکوستیک مربوط به هر جزء گاوسی در GMM،  $\mu_{iy}\mu_{ix}$ ، ساختار فرمانت بسیار مشابهی دارند. به دنبال این مشاهدات، یک تابع تاب‌خوردگی فرکانس خطی تکه‌ای  $W_i(f)$  با استفاده از موقعیت این شکل‌دهنده‌ها برای هر کلاس آکوستیک ایجاد می‌شود و تابع تاب‌خوردگی فرکانس برای قاب کامل شامل  $m$  کلاس‌های صوتی به دست می‌آید.

متعلق به من گاوسی است که توسط

(11)

$$W(f) = \sum_{i=1}^m w_i(x) W_i(f)$$

با  $W(f)$  تعیین شده، با فرض اینکه  $A(f)$  و  $\theta(f)$  اندازه و فاز تخمین‌زننده طیف قاب فعلی هستند، تغییرات طیفی را می‌توان با تاب برداشتن پوشش‌های منبع به شرح زیر انجام داد.

(12)

$$A_w(f) = A(W^{-1}(f)), \theta_w(f) = \theta(W^{-1}(f))$$

حتی پس از به دست آوردن طیف تاب‌خورده جریان فریم  $S_w(f)$  مانند بالا، توزیع انرژی همچنان با صدای هدف واقعی متفاوت است، زیرا شدت، پهنای باند و شیب طیفی تقریباً بدون تغییر باقی می‌مانند. تابع تبدیل (9) GMM در اینجا برای به دست آوردن استفاده می‌شود نسخه جدید طیف هدف  $S_g(f)$  از بردار LSF تبدیل شده  $F(x)$  طیف تبدیل نهایی برای فریم فعلی به دست می‌آید

(13)

$$S'(f) = G(f) S_w(f)$$

که در آن فیلتر تصحیح انرژی  $G(f)$  توسط

$$G(f) = \left| \frac{S_g(f)}{S_w(f)} \right| * B(f) \quad (14)$$

که از تابع پنجره هموارسازی  $B(f)$  برای انجام کانولوشن (عملگر  $*$ ) استفاده می کند. تابع  $B(f)$  می تواند تعادل بین شباهت و کیفیت گفتار تبدیل شده را با تنظیم شکل آن کنترل کند. تابع صاف کردن مثلی معمولاً در عمل استفاده می شود. [26]

برای تغییر بزرگی شکل دهنده ها، می توان پوشش تمام قطبی  $\hat{y} = F(x)$  تبدیل شده را به دست آورد و انرژی در باندهای خاص مورد علاقه را اندازه گیری کرد. انرژی قاب گفتار تبدیل شده در هر باند به سادگی با عوامل ضرب ثابت تصحیح می شود.

یک ویژگی مهم عروسی، فرکانس بنیادی  $f_0$ ، با تبدیل جهانی ساده اما محبوب تغییر یافته است [35]. از آنجایی که  $f_0$  از توزیع  $\log$ - نرمال پیروی می کند، میانگین  $\mu_{f_0}$  و واریانس  $\sigma_{f_0}^2 \log f_0$  را می توان در طول آموزش محاسبه کرد. سپس،  $f_0$  سیگنال گفتار تبدیل شده را می توان به صورت خطی با مقیاس بندی کرد

(15)

$$\log f_0^{(c)} = \mu_{f_0}^{(t)} + \frac{\sigma_{f_0}^{(t)}}{\sigma_{f_0}^{(s)}} \left( \log f_0^{(s)} - \mu_{f_0}^{(s)} \right)$$

که در آن حروف بالا (s) و (t) نشان دهنده منبع و هدف به ترتیب، و  $f(c)$  فرکانس اساسی است گفتار تبدیل شده در رابطه با مولفه تصادفی در رابطه (1)، آن است تبدیل شناخته شده است که به اندازه تبدیل هارمونیک / قطعی مرتبط نیست [13]. با این وجود، بهتر است مولفه تصادفی را برای گوینده هدف با استفاده از پارامترهای LSF دستگاه صوتی در فریم های صوتی، همانطور که در [36] شرح داده شده، پیش بینی کنیم. در این مرحله، ما شرح تئوری اصلی تبدیل صدا را برای سیستم خود به پایان رساندیم. در بخش بعدی، پیاده سازی بسیار کارآمد چنین الگوریتم هایی را شرح خواهیم داد.

## پیاده سازی کارآمد الگوریتم های اصلی

در مقایسه با لپ تاپ ها و رایانه های شخصی (رایانه های رومیزی شخصی)، دستگاه های تلفن همراه، به عنوان مثال، تلفن های هوشمند، معمولاً با حافظه، قدرت و منابع محاسباتی کاملاً محدود مشخص می شوند. با این وجود، هدف اصلی مطالعه ما توسعه یک سیستم تبدیل صدا کامل بر روی تلفن های همراه بدون پشتیبانی سمت سرور است. بنابراین، چالش کلیدی اجرای الگوریتم تبدیل به اندازه کافی کارآمد است به طوری که زمان مورد نیاز برای تبدیل صوتی، به ویژه فاز تبدیل، برای استفاده روزانه قابل قبول باشد. برای رفع این مشکل باید دو نکته زیر را عمده تاً در نظر گرفت: کارایی الگوریتم انتظار می رود روش هایی که برای مدل سازی، سنتز و تبدیل گفتار انتخاب می شوند، در زمان خود صرفه جویی کنند. در بخش دوم، ما

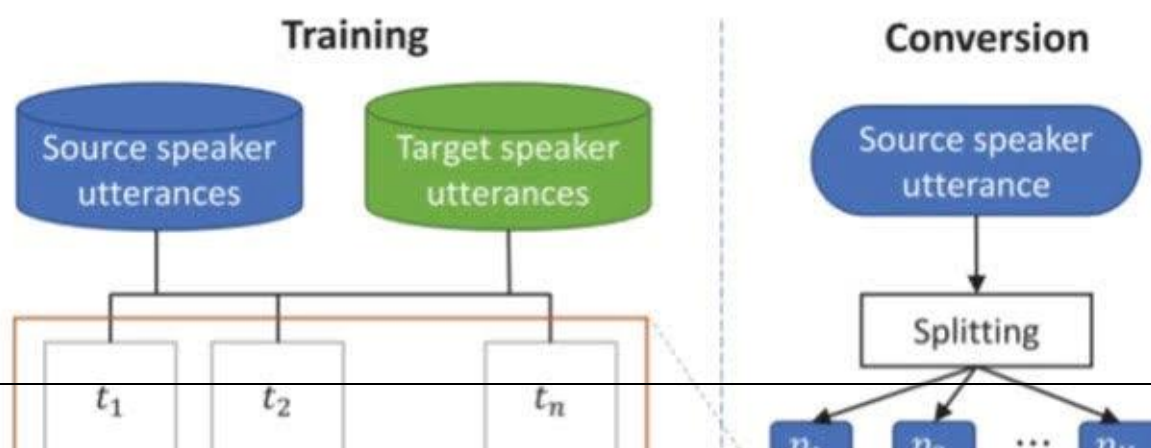
الگوریتم‌های درگیر برای سیستم تبدیل صدای خود را معرفی کرده‌ایم که تا حدی به دلیل بار محاسباتی کم و در عین حال عملکرد قابل قبول آن‌ها پذیرفته شده‌اند. به عنوان مثال، مدل قطعی به علاوه تصادفی که قبلاً برای مدل‌سازی گفتار توضیح داده شد، در واقع ناهمزمان است، که می‌تواند تحلیل را تا حد زیادی ساده‌تر کند، زیرا جداسازی دقیق دوره‌های سیگنال ضروری نیست. مهم‌تر از آن، این الگوریتم تبدیل می‌تواند حتی با تعداد کمی از نمونه‌های آموزشی، به عنوان مثال، ده‌ها جمله، به خوبی کار کند.

• کارایی پیاده‌سازی. این وظیفه‌ای است که ما در طول این مطالعه بیشتر بر آن تأکید می‌کنیم. به طور کلی، سخت‌افزار همیشه سریع‌تر از نرم‌افزار است. بنابراین، استفاده کامل از ویژگی‌های سخت‌افزاری پیشرفته موجود در تلفن‌های همراه رایج، از جمله پشتیبانی از محاسبات برداری و محاسبات موازی چند هسته‌ای، حیاتی است. علاوه بر این، انتخاب یک زبان برنامه‌نویسی مناسب و دستکاری هوشمند ماتریس‌ها نیز نقش مهمی در شتاب بازی می‌کند.

در این بخش، جنبه‌های اجرایی کلیدی در رابطه با کارایی محاسباتی را به تفصیل شرح می‌دهیم. امروزه، ارائه یک CPU چند هسته‌ای، حتی برای تلفن‌های هوشمند ارزان قیمت، برای تلفن‌های هوشمند تقریباً استاندارد است. به عنوان مثال، آیفون ۷ که در سپتامبر ۲۰۱۶ عرضه شد، از Apple A10 Fusion 64 بیتی سیستم روی تراشه (SoC) استفاده می‌کند که از دو هسته کم مصرف و دو هسته پرمصرف تشکیل شده است. به عنوان مثال دیگر، گوشی اندرویدی رده پایین Redmi 5 که در دسامبر ۲۰۱۷ توسط شیائومی عرضه شد، مجهز به تراشه اسنپدراگون ۴۵۰ با ۸ هسته است. بنابراین، برای استفاده کامل از ظرفیت محاسباتی تلفن‌های همراه، به‌ویژه پردازنده‌های چند هسته‌ای که به‌طور گسترده در تلفن‌های هوشمند در دسترس هستند، سیستم ما باید موازی‌سازی شود تا روی چندین هسته برای عملکرد تبدیل در زمان واقعی کار کند.

محاسبات موازی نوعی محاسبات است که در آن بسیاری از محاسبات یا اجرای فرآیندها به طور همزمان بر روی چندین هسته انجام می‌شود. به عبارت ساده، یک مسئله پیچیده بزرگ را می‌توان به مسائل کوچکتر و ساده‌تر تجزیه کرد، که می‌توان آن‌ها را به هسته‌های مختلف اختصاص داد تا همزمان حل شوند و در نتیجه کل زمان محاسبات مورد نیاز را کاهش دهند [۴۴]. با توجه به محدودیت‌های فیزیکی که از مقیاس‌بندی فرکانس جلوگیری می‌کند، محاسبات موازی مورد توجه گسترده‌تری قرار گرفته و به پارادایم غالب در معماری کامپیوتر تبدیل شده است، همانطور که با محبوبیت پردازنده‌های چند هسته‌ای نشان داده شده است [۴۵]. برای مثال، پردازنده‌های چند هسته‌ای برای رمزگشایی کارآمد ویدیوها مورد سوء استفاده قرار گرفته‌اند [۴۶]. در سیستم تبدیل صدای ما، محاسبات عمدتاً در دو نقطه موازی می‌شوند، مرحله مدل‌سازی و تحلیل گفتار در طول آموزش و تبدیل، که در شکل ۲ نشان داده شده است.

موازی‌سازی در طول مرحله آموزش، همانطور که در سمت چپ شکل ۲ نشان داده شده است، واضح و ساده است.



تصویر ۲: محاسبات موازی در مراحل آموزش و تبدیل

از آنجایی که تحلیل گفتار و ساخت ویژگی هر نمونه گفته مستقل است. بنابراین، می‌توانیم کل ۲ میلیون نمونه را از بلندگوی منبع و بلندگوی هدف بر روی هسته‌های C با ایجاد یک رشته ti بر روی هر هسته توزیع کنیم. به این ترتیب، در حالت ایده‌آل، نمونه‌های صوتی C می‌توانند همزمان بدون تداخل یکدیگر پردازش شوند. بنابراین، کل زمان اجرا بسیار کاهش می‌یابد زیرا زمان صرف شده توسط تجزیه و تحلیل گفتار و مهندسی ویژگی بخش بزرگی از کل هزینه زمان را اشغال می‌کند.

موازی کردن فاز تبدیل ظریف تر است. میانگین مدت جملات در مجموعه تقریباً ۴ یا ۵ ثانیه است، مشابه جملاتی که ما هر روز صحبت می‌کنیم. با توجه به ورودی گفته از بلندگوی منبع که باید تبدیل شود، ابتدا آن را به گروهی از پارتیشن‌های پیوسته  $2^i, i = 1, \dots, K$  تقسیم می‌کنیم، که طول آن‌ها تقریباً برابر است، مثلاً، در حدود tp ثانیه. سپس، مشابه مرحله آموزش، دوباره این پارتیشن‌های کوتاه را به هسته‌های مختلف برای پردازش همزمان اختصاص می‌دهیم، همانطور که در سمت راست در شکل ۲ نشان داده شده است. برای جبران مصنوعات احتمالی معرفی شده در مرز قطعه، همپوشانی اضافی وجود دارد. منطقه ای از مدت زمان تا ثانیه در هر طرف نقطه تقسیم بندی، نشان داده شده در شکل ۳. (a) پس از هر پارتیشن پی به

pic، ما این بخش‌ها را با استفاده از یک تابع لجستیک در یک گفتار تبدیل شده کامل ادغام می‌کنیم تا به یک انتقال صاف در اطراف نقطه پارتیشن دست یابیم. تابع لجستیکی که برای انتقال هموار در اینجا استفاده می‌شود توسط داده می‌شود

(16)

$$\varphi(x) = \frac{1}{1 + e^{-kx}}$$

جایی که k شیب منحنی را تنظیم می‌کند. در اجرای فعلی ما، طول همپوشانی است 40 میلی‌ثانیه، یعنی حدود ۶۴۰ نمونه. با فرض اینکه دو ناحیه همپوشانی اطراف نقطه تقسیم را از ۶۴۰- تا ۶۴۰.

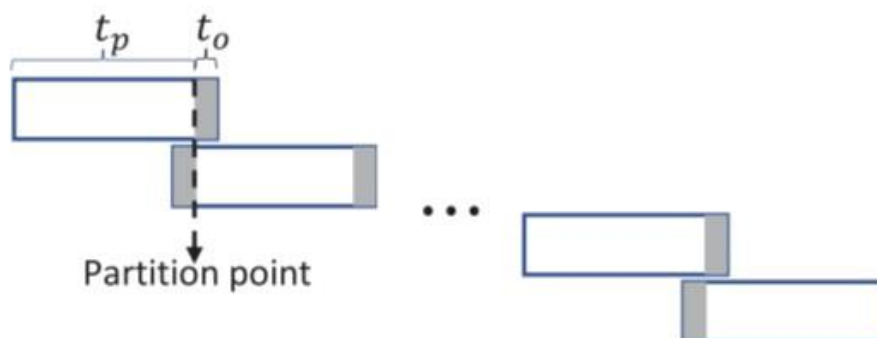
ایندکس کنیم و  $k = 0.015$  را در تابع لجستیک (۱۶) که در شکل ۳ (b) نشان داده شده است (انتخاب کنیم، نمونه در شاخص i دو نمونه متناظر را ادغام می‌کند.  $S_i^l$  در پارتیشن سمت چپ و آقا در پارتیشن سمت راست، توسط

(17)

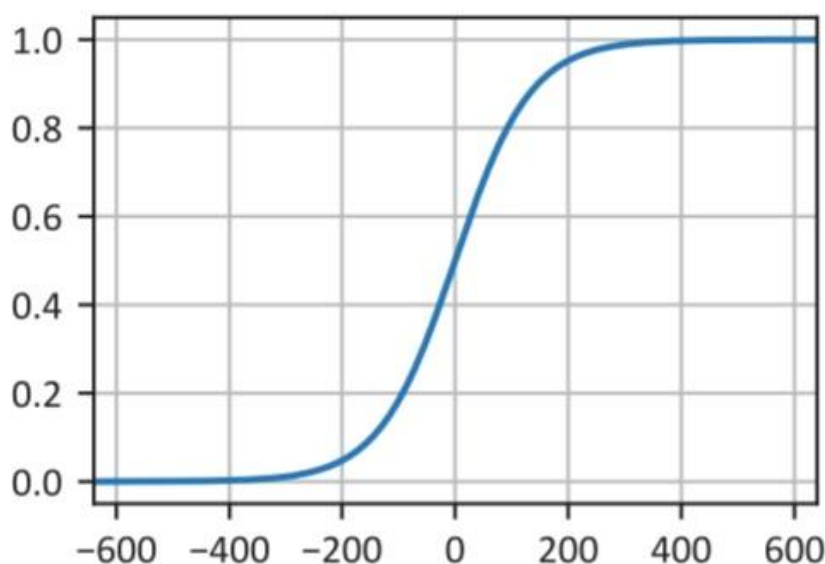
$$S_i = (1 - \varphi(i))s_i^l + \varphi(i)s_i^r$$

اگرچه در اجرای فعلی ما عمدتاً مراحل تحلیل گفتار، ساخت ویژگی و تراز فریم را به دلیل مصرف زیاد آنها

موازی می‌کنیم، باید توجه داشت که موازی‌سازی بیشتر همچنان می‌تواند در مراحل بعدی مانند تخمین پارامترهای GMM اعمال شود. معادله (۸). به عنوان مثال، در کار [47] Wojciech Kwedlo، یک موازی‌سازی حافظه مشترک از الگوریتم استاندارد EM بر اساس تجزیه داده‌ها برای یادگیری پارامترهای GMM در یک سیستم چند هسته‌ای برای عملکرد بالاتر پیشنهاد شده است. بنابراین، ما در حال برنامه‌ریزی برای پیاده‌سازی چنین ویژگی‌هایی برای تسریع بیشتر سیستم تبدیل صدای خود در نسخه بعدی آن هستیم.



Split a sentence for parallel conversion  
(a)



Logistic function for segment merging  
(b)

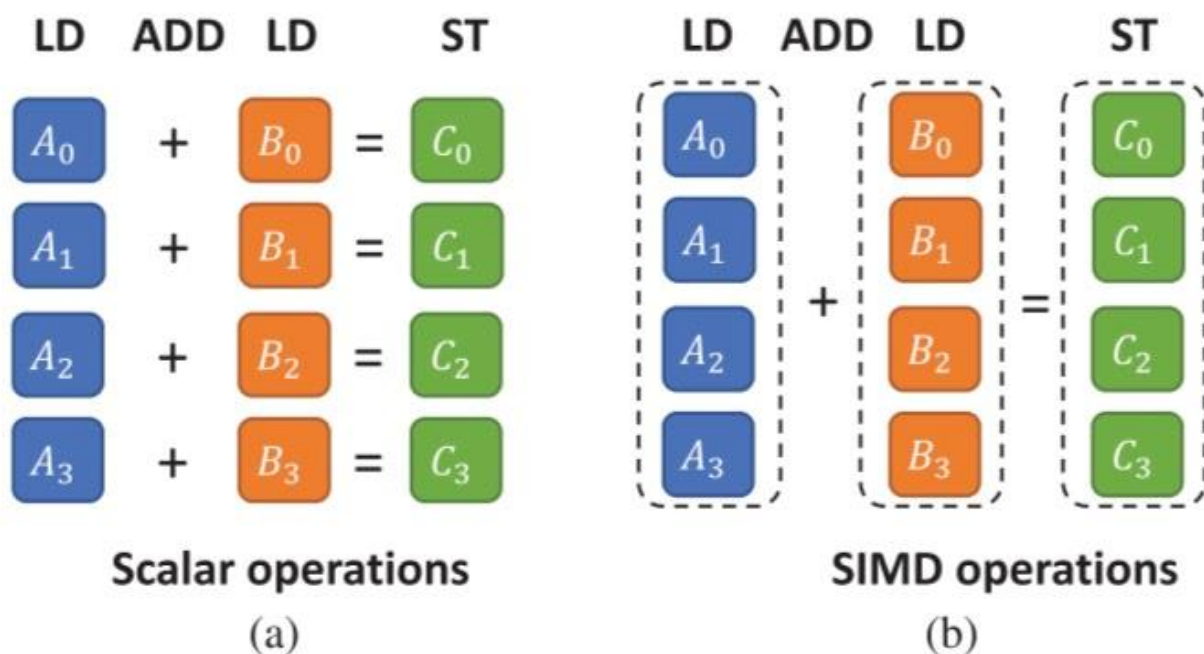
تصویر ۳: موازی‌سازی چند هسته‌ای در مرحله تبدیل. اول، یک جمله ورودی به چند بخش تقسیم می‌شود و آن بخش‌ها به طور همزمان تبدیل می‌شوند. سپس، بخش‌های تبدیل شده با استفاده از یک تابع لجستیک برای مجموع وزنی به آرامی ادغام می‌شوند. (الف) یک جمله را برای تبدیل موازی تقسیم کنید (ب) تابع لجستیک برای ادغام بخش.

## برداری از طریق SIMD

پس از اینکه کار محاسباتی روی چندین هسته گسترده شد، هدف بعدی ما این است که توان عملیاتی تک هسته‌ای را تا حد امکان بالا ببریم. مجدداً، به دلیل استفاده زیاد از ماتریس‌ها در الگوریتم‌های زیربنایی، ما



باید به کاوش و استفاده کامل از ظرفیت سخت افزار، به ویژه موازی سازی درون هسته ای برای عملیات روی آرایه ها (بردارها) داده ها بپردازیم. در تلفن های همراه مدرن، موازی سازی درون هسته ای، یا بردارسازی، عمدتاً توسط SIMD پشتیبانی می شود که به طور گسترده در پردازنده های امروزی موجود است [44]، [48]. اکثر پردازنده های موبایل، از جمله پردازنده های آیفون و اندروید، با معماری ARM طراحی شده اند. به عنوان مثال، آیفون ۷ از سیستم روی تراشه ۶۴ بیتی Apple A10 Fusion استفاده می کند که مبتنی بر معماری ARMv8-A با پشتیبانی از دستورالعمل های ۶۴ بیتی است. به ویژه، پردازنده های مدرن موبایل ARM، مانند تمامی سری های Cortex-A8، معمولاً از افزونه های پیشرفته SIMD، معروف به NEON، پشتیبانی می کنند که یک مجموعه دستورات SIMD ترکیبی ۶۴ و ۱۲۸ بیتی است که شتاب استاندارد را برای پردازش رسانه و سیگنال ارائه می کند. برنامه های کاربردی NEON [48] می تواند از اعداد صحیح ۸، ۱۶، ۳۲ و ۶۴ بیتی پشتیبانی کند.



تصویر ۴: عملیات اسکالر در مقابل SIMD برای افزودن های متعدد.

و داده های ممیز شناور تک دقیق (۳۲ بیتی) و عملیات SIMD تا ۱۶ عملیات را همزمان فراهم می کند. در یک مطالعه بر روی تثبیت کننده تصویر دیجیتال برای دستگاه های تلفن همراه، محاسبات بردار حرکت و محاسبات FFT با استفاده از موتور نئون SIMD موجود در CPU ARM تسریع می شوند و بردار حرکت جهانی برای هر فریم را می توان در کمتر از ۲۰ میلی ثانیه محاسبه کرد [۴۹]. در معماری های x86/64 که معمولاً برای رایانه های شخصی استفاده می شود، پردازنده ها معمولاً مجموعه دستورالعمل های SIMD مبتنی بر SSE یا AVX را ارائه می کنند.

به طور خلاصه، SIMD عملیات مشابهی را روی چندین عنصر داده از یک نوع که در یک بردار تراز شده اند به طور همزمان انجام می دهد. از نظر تئوری، اگر بتوانیم عناصر داده  $q$  را به طور کلی با یک SIMD پردازش کنیم، آنگاه سرعت افزایش در مقایسه با پردازش متوالی از طریق دستورالعمل های اسکالر به  $q$

نزدیک است، اگرچه سرعت واقعی به طور بالقوه توسط پهنای باند حافظه محدود می شود. شکل ۴ یک عملیات رایج در پردازش چند رسانه ای را نشان می دهد، که در آن مقدار یکسان به تعداد زیادی از عناصر داده اضافه می شود، که با دستورالعمل های اسکالر و دستورالعمل های SIMD به طور جداگانه برای مقاصد مقایسه اجرا می شوند. در زمینه پردازش صدا، داده های صوتی معمولاً در انواع اعداد صحیح ۱۶ بیتی ارائه می شوند. با فرض اینکه از رجیسترهای NEON 64 بیتی برای برداری استفاده می کنیم، می توانیم چهار عدد صحیح ۱۶ بیتی را همزمان در این ثبات واحد بسته بندی کنیم، همانطور که در شکل ۴ (b) نشان داده شده است. فرآیندی که فقط با استفاده از دستورالعمل های اسکالر در شکل ۴ (a) نشان داده شده است. برای انجام یک جمع، ابتدا دو عملیات بارگذاری (LD) برای خواندن دو عدد A و B از حافظه اجرا می شود. سپس، یک افزودنی (ADD) برای به دست آوردن نتیجه C دستور داده می شود. در نهایت، مجموع C با یک دستور ذخیره (ST) به حافظه بازگردانده می شود. در مقابل، دستورالعمل های بارگذاری، افزودن و ذخیره سازی تقویت شده SIMD می توانند روی چهار عدد صحیح به طور همزمان اجرا شوند که در شکل ۴ (b) مشخص شده است. بنابراین، در مجموع، ۱۶ دستورالعمل اسکالر از طریق برداری SIMD به تنها چهار دستورالعمل کاهش می یابد که سرعت تئوری چهار را به دست می دهد.

برای استفاده از SIMD در برنامه نویسی عملی و توسعه نرم افزار، عمدتاً سه راه وجود دارد: (۱) کدگذاری در دستورالعمل های اسمبلی سطح پایین به طور مستقیم. (۲) استفاده از توابع ذاتی با رابط های C که کدهای اسمبلی را به صورت داخلی بسته بندی می کنند. و (۳) استفاده از بردار سازی خودکار در صورتی که توسط کامپایلرها پشتیبانی شود [۴۸]. به طور کلی، برنامه نویسی در زبان اسمبلی به طور مستقیم، علیرغم عملکرد بسیار زیاد، کار فشرده و بسیار مستعد خطا است. از سوی دیگر، بردار سازی خودکار می تواند به طور کامل توسط خود کامپایلر بدون نیاز به مداخله انسانی انجام شود. با این حال، در حال حاضر حتی کامپایلرهای پیشرفته هستند

به اندازه کافی هوشمند نیست، و در نتیجه، تنها بخش کوچکی از کدها می توانند به طور خودکار بردار شوند. بنابراین، مطالعه ما بر عملکردهای ذاتی به عنوان یک مبادله بین عملکرد سیستم و کار دستی متکی است. شایان ذکر است که در یک پیاده سازی واقعی، وابستگی به توابع ذاتی لزوماً به این معنا نیست که ما باید آنها را مستقیماً توسط خودمان در برنامه نویسی فراخوانی کنیم. قابل توجه است که بخش عمده ای از الگوریتم های تجزیه و تحلیل گفتار و تبدیل در محاسبات ماتریسی فرموله می شوند. بنابراین، ما می توانیم به برخی از کتابخانه های جبر خطی بالغ متوسل شویم که رابط های سطح بالا را برای تسهیل توسعه برنامه ها در معرض دید قرار می دهند و در عین حال از توابع ذاتی برای بردار کردن محاسبات ماتریس/آرایه تا حد امکان استفاده می کنند. یعنی استفاده داخلی از دستورالعمل های SIMD برای کاربران تقریباً شفاف است. جزئیات مربوطه در زیر بخش بعدی توضیح داده شده است

در دو بخش فرعی بالا، موازی سازی سطح هسته و سطح دستورالعمل را برای تسریع الگوریتم های اصلی برای تبدیل صدا توضیح داده ایم. هدف از این مطالعه توسعه یک سیستم تبدیل صدا کارآمد در دستگاه های تلفن همراه است، نه محدود به یک نوع دستگاه خاص مانند گوشی های iPhone، iPad یا Android. اگرچه کیت های توسعه رابط کاربری گرافیکی در سیستم های عملیاتی مختلف تلفن همراه، مانند iOS برای آیفون و اندروید برای اکثر تلفن های هوشمند دیگر، تفاوت زیادی با یکدیگر دارند، اما ما می خواهیم اجرای

الگوریتم اصلی این سیستم مستقل از پلتفرم باشد به طوری که این امر بسیار مهم باشد. بخش را می توان به راحتی در چندین پلت فرم با حداقل تغییرات منتقل کرد. از این رو، ما عملکرد اصلی و رابط کاربری گرافیکی را در معماری سیستم خود جدا می کنیم. برای برآوردن نیازهای بی طرفی پلت فرم و بازده بالا، زبان برنامه نویسی **C++** بهترین گزینه برای کدگذاری الگوریتم های اصلی است. به عنوان یک زبان همه منظوره، **C++** عملکرد، کارایی و انعطاف پذیری را برجسته می کند، یعنی امکان دستکاری حافظه در سطح پایین را فراهم می کند و در عین حال انتزاعات سطح بالا را به شیوه ای شی گرا ارائه می کند. در زیر به طور خلاصه نحوه دستیابی به موازی سازی چند هسته ای و مبتنی بر **SIMD** در **C++** را شرح می دهیم.

برای اعمال موازی سازی مبتنی بر رشته برای محاسبات چند هسته ای، یعنی برنامه نویسی چند رشته ای در **C++**، ساده ترین راه استفاده از **OpenMP API** استاندارد است که از برنامه نویسی موازی با حافظه مشترک چند پلتفرمی در **C** و **C** پشتیبانی می کند [50]. **C++** تنها با استفاده از دستورالعمل های ساده. با این حال، **OpenMP** هنوز در **iOS**، سیستم عامل تلفن همراه آیفون یا آی پد، پشتیبانی ضعیفی دارد. در عوض، همانطور که اپل توصیه می کند، باید از فناوری **Grand Central Dispatch (GCD)** برای موازی سازی وظایف استفاده کنیم، که به طور ویژه برای **iOS** طراحی شده و بسیار بهینه شده است. به طور خلاصه، آن چارچوب های موازی ساده هنوز تا حدودی به پلت فرم وابسته هستند. خوشبختانه، این مشکل قابل حمل را می توان با استاندارد جدید **C++ 11** که در سال ۲۰۱۱ منتشر شد، حل کرد که با معرفی یک کتابخانه رشته جدید، پشتیبانی از برنامه نویسی چند رشته ای را اضافه می کند. به عبارت ساده، ما می توانیم با نمونه سازی کلاس **std::thread** با ارسال تابعی که نشان دهنده کارهایی است که باید در این رشته انجام شود، یک رشته جدید ایجاد کنیم. باید توجه ویژه ای شود

اجتناب از مسابقه داده در چنین برنامه های چند رشته ای با استفاده از اصول اولیه همگام سازی مانند و برای محافظت از داده های مشترک. خوانندگان علاقه مند می توانند برای جزئیات بیشتر در مورد برنامه نویسی چند رشته ای در **C++** به این کتاب عالی [۵۱] مراجعه کنند.

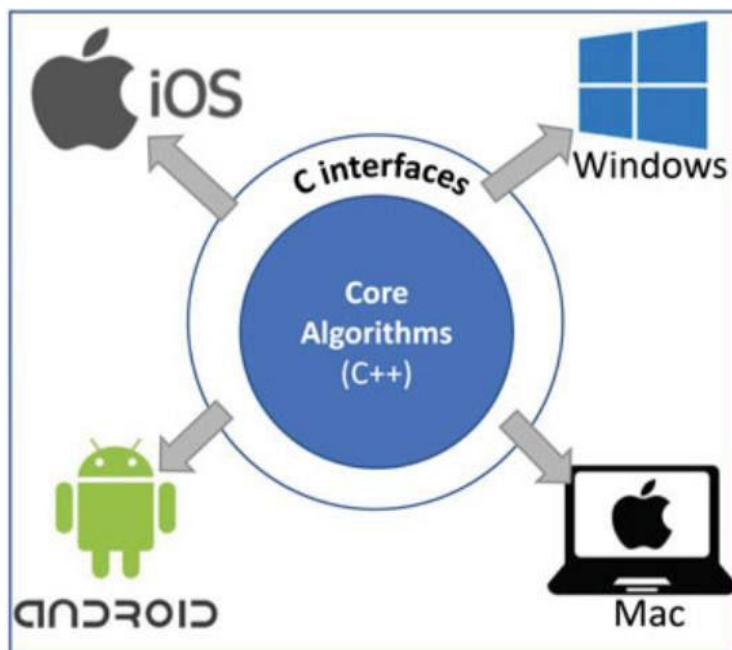
در مورد برداری بر روی هر هسته با استفاده از دستورالعمل های **SIMD**، ما به کتابخانه معروف **Eigen**، یک کتابخانه الگوی **C++** با کارایی بالا برای جبر خطی، که در بسیاری از پروژه های صنعتی استفاده شده است، تکیه می کنیم [۵۲]. ما **Eigen** را انتخاب می کنیم زیرا از تمام اندازه های ماتریس پشتیبانی می کند و الگوریتم های رایج برای تجزیه ماتریس، حل کننده های معادلات عددی و سایر الگوریتم های مرتبط مانند تبدیل فوریه سریع (**FFT**) را ارائه می دهد. مهمتر از آن، **Eigen** با انجام بردار سازی صریح داخلی برای مجموعه های دستورات **SSE**، **AVX** و **ARM NEON** به اندازه کافی سریع است و این پیچیدگی را از کاربران پنهان می کند. به این معنا که ما می توانیم با اجرای ساده الگوریتم ها با استفاده از ماتریس ها و عملیات جبر خطی ارائه شده توسط **Eigen**، تا حد امکان به موازی سازی در سطح دستورالعمل دست یابیم. به عنوان یک مزیت جانبی، **Eigen** رابط هایی مشابه **MATLAB** ارائه می دهد که می تواند برنامه نویسی را ساده کرده و کدها را به ویژه برای مهندسان خواناتر کند.

اکنون می توانیم معماری نرم افزاری سیستم تبدیل صدای خود را که در شکل ۵ نشان داده شده است، نمایانگر داشته باشیم. واضح است که سیستم ما الگوریتم های اصلی (موتور محاسباتی) و رابط کاربری گرافیکی (**GUI**) را به **maxi** تقسیم می کند. - قابلیت حمل پلت فرم **mize** در حالت ایده آل، هیچ تغییری در کد

منبع الگوریتم‌های هسته‌ای هنگام انتقال به پلتفرم دیگر مورد نیاز نیست، و ما فقط باید کدهای **C++** را با کامپایلرهای خاص پلتفرم برای استقرار مجدد کامپایل کنیم. برعکس، به طور کلی، بخش رابط کاربری گرافیکی با پلتفرم‌های خاصی مانند **iOS**، اندروید و ویندوز به شدت مرتبط است و نمی‌توان آن را منتقل کرد. به عنوان مثال، در **iOS**، ما رابط کاربری گرافیکی را با فریم ورک **UIKit** با استفاده از زبان **Swift** یا **Objective-C** می‌سازیم، در حالی که اندروید اجزای رابط کاربری خود را فراهم می‌کند که با زبان جاوا یا کاتلین قابل دسترسی هستند. یکی از مشکلات هسته محاسباتی مبتنی بر **C++** این است که **C++** فقط با برخی از زبان‌ها از جمله **Swift**، زبان اصلی توسعه **iOS**، قابلیت همکاری ضعیف یا حتی بدون همکاری را پشتیبانی می‌کند. بنابراین، ما تصمیم می‌گیریم رابط‌های ورودی/خروجی الگوریتم‌های اصلی را با استفاده از زبان برنامه‌نویسی **C** پیچیم، که می‌تواند به راحتی با تمام زبان‌های رایج تعامل داشته باشد. در برنامه کامل تلفن همراه، رابط کاربری گرافیکی با انتقال اطلاعات عمدتاً از طریق قطار و تبدیل توابع ارائه شده به عنوان رابط‌های **C** با الگوریتم‌های اصلی ارتباط برقرار می‌کند. به این ترتیب، موتور محاسباتی، یعنی الگوریتم‌های هسته، می‌تواند به خوبی روی هر چهار سیستم عملیاتی نشان داده شده در شکل ۵ همانطور که ما آزمایش کرده‌ایم، دو مورد برای دستگاه‌های تلفن همراه و دو مورد برای رایانه‌های شخصی، کار کنند، بنابراین هدف ما از استقلال پلت فرم را محقق می‌کند.

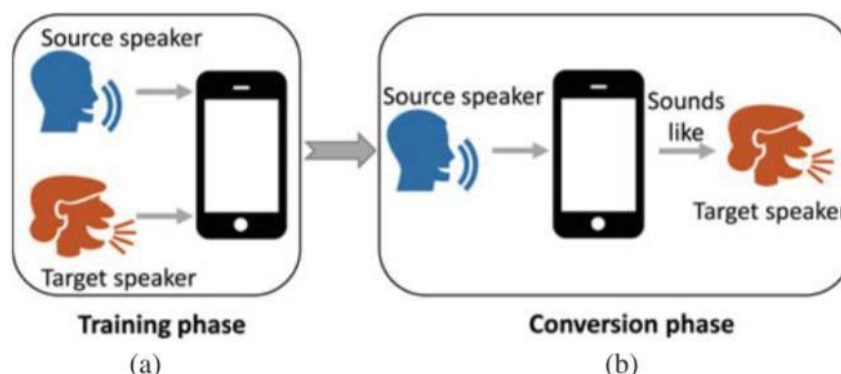
## توسعه برنامه های iOS

پس از الگوریتم‌های اصلی سیستم تبدیل صدا با **C++** پیاده سازی می شوند، کار باقی مانده است یک رابط کاربری گرافیکی دوستانه در تلفن های همراه ایجاد کنید، یعنی توسعه دهید.



تصویر ۵: مروری بر معماری نرم افزار الگوریتم‌های اصلی (موتور محاسباتی) برای قابلیت حمل در **C++** پیاده سازی شده اند. رابط گرافیکی کاربر (GUI) در بالای موتور قرار دارد و ممکن است با زبان‌ها/کتابخانه

های مختلف در پلتفرم های مختلف ساخته شود. این موتور رابط های C را برای استفاده توسط رابط کاربری گرافیکی در چهار سیستم عامل رایج iOS، Android، Windows و Mac در معرض دید قرار می دهد.

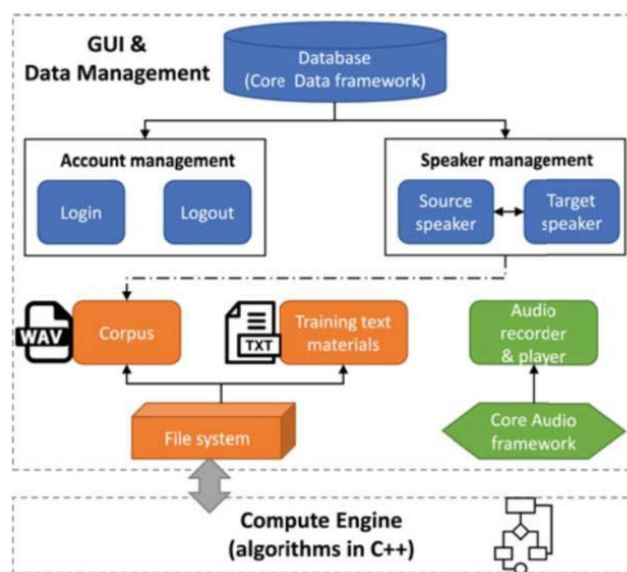


تصویر ۶: گردش کار برنامه تبدیل صدا در تلفن های همراه

اپلیکیشن موبایل، برای تسهیل دسترسی به این سامانه. همانطور که قبلاً در شکل ۵ نشان داده ایم، می توانیم رابط کاربری گرافیکی را با زبان ها و ابزارهای مختلف بر روی پلتفرم های مختلف توسعه دهیم که به خوبی از موتور محاسباتی در معماری ما جدا شده است. در این بخش، به عنوان یک مثال خاص، ما توسعه یک برنامه (برنامه) iOS را ارائه می کنیم که بر روی آیفون ۷، احتمالاً محبوب ترین تلفن هوشمند در سال ۲۰۱۷، مستقر شده است. به طور کلی، به عنوان یک نمونه اولیه برای تأیید امکان سنجی صوتی نسخه بر روی تلفن های همراه، برنامه ما بر قابلیت استفاده و مختصر تمرکز دارد. اساسی ترین استفاده از این برنامه فقط شامل دو مرحله است که در شکل ۶ نشان داده شده است. اولاً، هم گوینده منبع و هم بلندگوی هدف تعداد کمی از جملات را می خوانند در حالی که گفته های آنها توسط تلفن ضبط می شود تا یک پیکره موازی برای موارد بعدی تشکیل شود. آموزش مدل تبدیل ثانیاً، بلندگوی منبع می تواند دوباره چیزی را با برنامه صحبت کند و سعی می کند این پیام را به گونه ای تبدیل کند که به نظر برسد که توسط بلندگوی مورد نظر صحبت می شود. در ادامه جزئیات بیشتری در مورد ساختار ماژولار و برخی مسائل طراحی خاص این اپلیکیشن iOS ارائه شده است.

## نمای کلی ماژول های کاربردی

در مهندسی نرم افزار، به خوبی شناخته شده است که ماژولار -ization یک اصل اساسی برای نرم افزارهای مقیاس بزرگ است



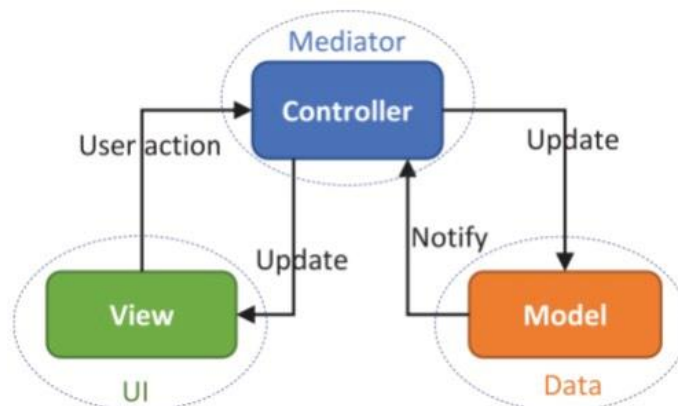
تصویر ۷: طرح کلی از ماژول ها در برنامه کاربردی

توسعه. ماژول های اصلی در سیستم ما در شکل ۷ نشان داده شده اند. برای مدیریت حساب های کاربری، از جمله ثبت نام، ورود به سیستم و خروج، اطلاعات کاربر را در یک پایگاه داده موبایل سبک مانند SQLite ذخیره می کنیم که دسترسی به آن واسطه است. توسط چارچوب iOS Core Data برای ساده کردن کدنویسی. به طور مشابه، برای یک بلندگوی منبع (یعنی کاربر)، چندین بلندگوی هدف را می توان آموزش داد، و ارتباط بین بلندگوی منبع و بلندگوی هدف نیز در پایگاه داده ثبت می شود. بیشتر فضای ذخیره سازی هارد دیسک این برنامه توسط مجموعه ساخته شده برای آموزش تبدیل صوتی اشغال شده است. ما مطالب متنی را ارائه می دهیم که هر کدام حاوی حدود ۲۰ جمله عادی است تا توسط یک جفت گوینده منبع و گوینده هدف خوانده شود. در حالی که یک سخنران در حال خواندن متن است، صدای او می تواند به راحتی توسط برنامه ما ضبط شود (شکل ۶) و با پشتیبانی از چارچوب iOS Core Audio به عنوان فایل های WAV تک صدایی با فرکانس نمونه برداری ۱۶ کیلوهرتز ذخیره شود.

در نهایت، باید از شکل ۷ متوجه شد که در حال حاضر در سیستم ما تعامل بین لایه رابط کاربری گرافیکی و موتور محاسباتی، یعنی الگوریتم های اصلی پیاده سازی شده در C++، کاملاً واضح و مختصر است. در طول آموزش، لایه رابط کاربری گرافیکی فقط باید به موتور بگوید فایل های WAV که مجموعه آموزشی را تشکیل می دهند چیست و موتور یک فایل مدل مطابق با تابع تبدیل تولید می کند. به طور مشابه، در نقطه تبدیل، لایه رابط کاربری گرافیکی موتور را در مورد فایل صوتی از بلندگوی مبدأ که قرار است تبدیل شود، مطلع می کند و موتور نهایتاً گفتار تبدیل شده را در یک فایل WAV جدید ذخیره می کند و به لایه رابط کاربری گرافیکی اطلاع می دهد. از مسیر فایل ما می توانیم صدای تبدیل شده را به سادگی با پخش این فایل صوتی بشنویم. علاوه بر این، همچنین می توان تنظیم کرد که پس از اتمام به طور خودکار پخش شود.

## طراحی های خاص

ایجاد یک عملکرد خوب یک کار بسیار چالش برانگیز است برنامه iOS، که شامل طراحی UI، منطق تجاری است



تصویر ۸: الگوی طراحی MVC

مدل سازی، کدنویسی با Swift/Objective-C، اشکال زدایی XCode و بسیاری کارهای خسته کننده دیگر. در اینجا، جزئیات برای حفظ فضا حذف شده است. در ادامه، ما فقط دو نگرانی طراحی خاص در مورد توسعه اپلیکیشن را به عنوان نمایندگان کل فرآیند فهرست کرده و به آن می پردازیم.

1. الگوی طراحی MVC برای رابط کاربری در توسعه برنامه های کاربردی مبتنی بر رابط کاربری گرافیکی، بهترین روش شناخته شده جداسازی داده ها از ارائه آن است، زیرا ارائه ممکن است با توجه به نیازهای خاص، به عنوان مثال، در نمودار، بسیار متفاوت باشد. یا در یک جدول از زمان اولین عرضه iOS، اپل توصیه کرده است که الگوی طراحی Model-View-Controller (MVC) را برای تأکید بر این بهترین روش، که در شکل ۸ مشخص شده است، اتخاذ کند. view. نشان دهنده چیزی قابل مشاهده در رابط کاربری است که برای نمایش داده ها اختصاص داده شده است. با این حال، مدل و نما نمی توانند به طور مستقیم به منظور کاهش جفت شدن ارتباط برقرار کنند. در عوض، آنها توسط شی کنترلر، معمولاً از طریق الگوی نمایندگی، واسطه می شوند.

در عمل، همانطور که در شکل ۷ نشان داده شده است، یک شی مدل را برای دسترسی به داده های ذخیره شده در پایگاه داده یا سیستم فایل تعریف می کنیم. کنترل کننده سعی می کند نما را با توجه به داده های تغییر یافته به روز کند. در جهت دیگر، هنگامی که کاربر با نمای تعامل برقرار می کند، اقدام او به کنترلر منتقل می شود، که بیشتر مدل را در مورد قصد کاربر آگاه می کند، به عنوان مثال، به روز رسانی یا حذف داده ها. علیرغم سادگی، لایه های یک برنامه رابط کاربری گرافیکی معمولی را می توان با الگوی MVC جدا کرد تا سازماندهی بهتر و ترویج استفاده مجدد از کد را تشویق کند. در برنامه ما، هر پنجره با الگوی MVC که در بالا توضیح داده شد ساخته شده است.

2. ارائه داده محور در نمای مجموعه برای نمایش چندین اطلاعات صوتی و آیتیم های بلندگوی هدف، نمای جدول و مجموعه نماها به طور گسترده در رابط کاربری گرافیکی برنامه ما استفاده می شود. برای



انعطاف پذیری بهتر و قابلیت استفاده مجدد برنامه، بهتر است داده ها را از تجسم آن از طریق سبک -UI Data-Operation جدا کنید. در چارچوب توسعه iOS، داده ها در یک شی منبع داده ذخیره می شوند و عملیات توسط یک شی نماینده نشان داده می شود. توجه داشته باشید که نمای جدول و نمای مجموعه در iOS UIKit همگی با مکانیزم تخصصی MVC طراحی شده اند.

منبع داده تنها مسئول ارائه داده ها است و نمی داند داده ها چگونه نمایش داده می شوند. تفویض اختیار به اشیا فرصتی می دهد تا ظاهر و حالت خود را با تغییراتی که در جاهای دیگر برنامه اتفاق می افتد، که معمولاً توسط اقدامات کاربر ایجاد می شود، هماهنگ کنند. مهمتر از آن، تفویض این امکان را برای یک شی فراهم می کند تا رفتار شیء دیگر را بدون نیاز به ارث بردن از آن تغییر دهد. با این الگوی طراحی، می توانیم منبع داده و ارائه آن را بهتر جدا کنیم. هنگامی که داده ها در جایی به روز شوند، UI به طور خودکار پاسخ می دهد تا تغییرات داده را منعکس کند. از این رو، چنین سبک ارائه مبتنی بر داده نیز به طور گسترده در برنامه ما به کار گرفته شده است تا توسعه رابط کاربری پیچیده را تسهیل کند.

## آزمایش ها و نتایج

ما سیستم تبدیل صدا را به شیوه ای مناسب از پایین به بالا توسعه داده ایم زیرا دو لایه در شکل ۷ تقریباً کاملاً جدا شده اند. یعنی موتور محاسباتی مسئول محاسبات فشرده ابتدا با تمرکز بر دقت و کارایی با آزمایش و بهبود مکرر ساخته شد. گام بعدی طراحی رابط کاربری و مدیران داده مرتبط برای یک برنامه تلفن همراه با توجه به قابلیت استفاده و مختصر بودن بود. برنامه تلفن همراه تمام شده، به نام Voichap، با موفقیت بر روی یک دستگاه آیفون ۷ مستقر شد. در ادامه، ابتدا سرعت الگوریتم هسته فعال شده توسط موازی سازی مبتنی بر چند هسته و بردار را ارزیابی می کنیم. سپس، رابط کاربری این اپلیکیشن را شرح می دهیم و استفاده از آن را از طریق مثال های خاص نشان می دهیم. در نهایت، کارایی و اثربخشی این سیستم تبدیل صدای تلفن همراه را با اندازه گیری زمان اجرای آن و انجام تست کیفیت تبدیل امتیازدهی شده توسط ۱۰ شنونده، ارزیابی می کنیم.

الف) الگوریتم های هسته ای با سرعت بالا نتایج تست را افزایش می دهند

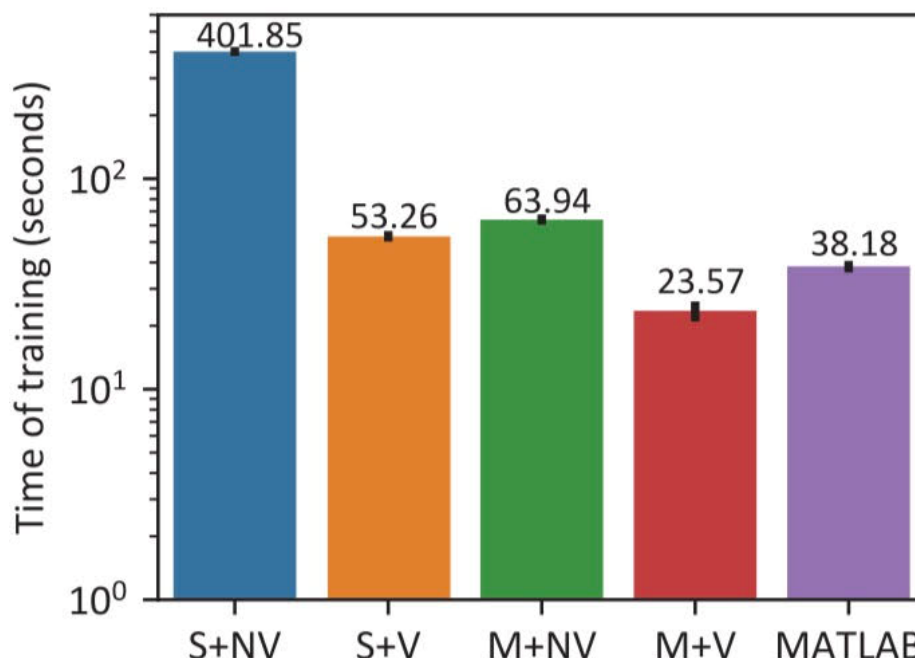
برای ارزیابی کارایی پیاده سازی خود با استفاده کامل از موازی سازی، عملکرد موتور را با تنظیمات محیطی مختلف آزمایش کردیم. برای سادگی، آزمایش ها را بر روی یک رایانه شخصی ویندوز ۱۰ با یک CPU Intel Core i7 با ۴ هسته سخت افزاری انجام دادیم، زیرا الگوریتم ها برای اولین بار با Visual Studio IDE قدرتمند در ویندوز ۱۰ توسعه و آزمایش شدند. این همچنین نشان می دهد که الگوریتم اصلی ما پیاده سازی در واقع پلتفرم آگنوستیک است، زیرا می توانیم همان کد را در iOS (iPhone 7) بدون هیچ تغییری اجرا کنیم. در هر آزمون اندازه گیری زمان اجرا ذکر شده در زیر، آزمایش ها برای تخمین دقیق تر ۱۰ بار تکرار شد.

### 1) عملکرد کلی

ما فرآیند آموزش را برای ۲۰ بیان زمان بندی کردیم که هر کدام به طور متوسط ۴,۵ ثانیه طول کشید، هم از بلندگوی منبع و هم از بلندگوی هدف، که پیکربندی پیش فرض برنامه تلفن همراه ما است. در این آزمون، تعداد مؤلفه های گاوسی مورد استفاده در مدل  $m = 4$  (7) انتخاب شده است. به عنوان خط پایه، ما همچنین



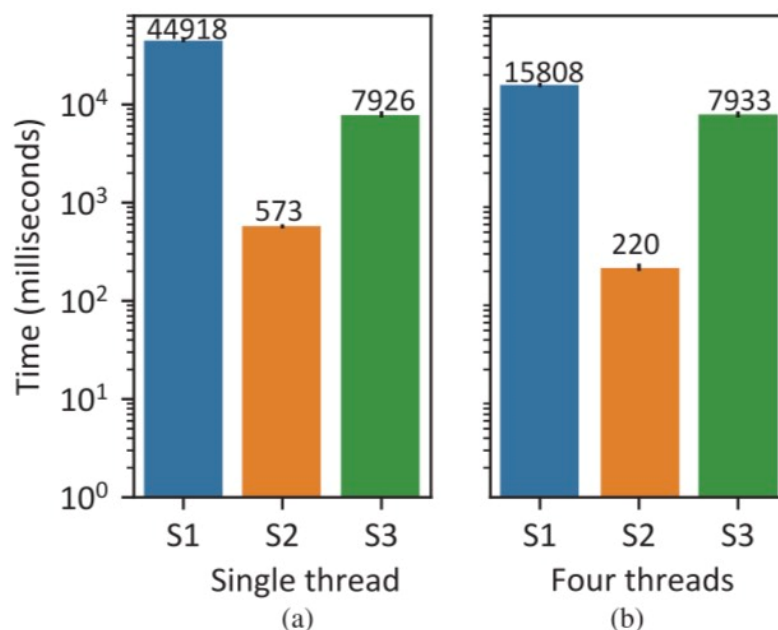
زمان اجرای کد اصلی متلب را در همان رایانه شخصی اندازه‌گیری کردیم. در اینجا لازم به یادآوری است که بسیاری از توابع داخلی متلب در واقع روال‌های C یا Fortran بسیار بهینه شده هستند که ممکن است به خوبی به صورت موازی در داخل پیاده‌سازی شوند. بنابراین، معمول است که یک برنامه عددی C++ ساده نتواند همتای خود را در MATLAB در بازه زمانی شکست دهد. آمار زمان اجرای اندازه‌گیری شده



در شکل ۹ نشان داده شده است. علاوه بر این، باید توجه داشته باشیم که روش تاب برداشتن فرکانس فقط فریم‌های صدادار را تغییر می‌دهد. بنابراین، تعداد فریم‌های صوتی در گفته‌ها ممکن است مورد توجه بیشتری باشد. در مجموعه آموزشی ما شامل ۲۰ گفتار، میانگین تعداد فریم‌های صوتی در هر گفته ۳۰۲,۶ و تغییر فریم ۸ میلی ثانیه (۱۲۸ نمونه) است. بنابراین، میانگین طول سیگنال صوتی در مجموعه آموزشی تنها ۲,۴۲ ثانیه است. به طور مشابه، می‌توانیم میانگین تعداد فریم‌های صدادار در مجموعه آزمایشی را برای تبدیل بشماریم، که در حدود ۲۹۷,۵ است که مربوط به میانگین مدت زمان ۲,۳۸ ثانیه است، اگر فقط بخش‌های صداگذاری شده را در نظر بگیریم برای پیاده‌سازی واقعا کارآمد در تلفن‌های همراه

**تصویر ۹:** زمان اجرای مرحله تمرین با تنظیمات مختلف (فاصله اطمینان ۹۵). (از چپ به راست، S+NV: C++ تک رشته‌ای بدون بردار، S+V: C++ تک رشته‌ای با برداری، M+NV: C++ چند رشته‌ای بدون برداری، M+V: C++ چند رشته‌ای با بردارization، -، MATLAB: 64 بیتی MATLAB 2016 با تنظیمات پیش فرض)

در اجرای فعلی مرحله آموزش، دو مرحله اول شامل تجزیه و تحلیل گفتار، ساخت ویژگی، و تراز قاب (به شکل ۱ مراجعه کنید) روی چند هسته موازی می‌شوند. در حالت ایده آل، هر گفته یا جفت منبع-هدف



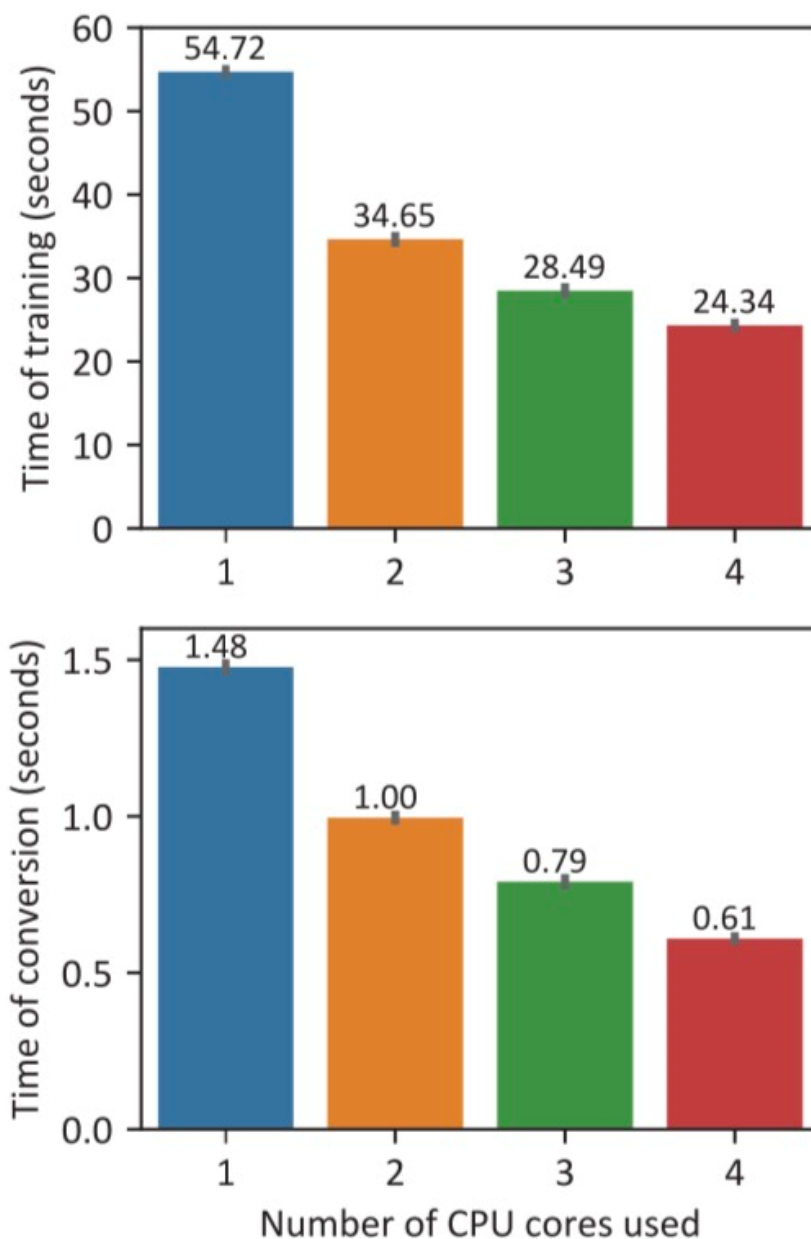
**تصویر ۱۰:** زمان اجرای هر مرحله در طول مرحله تمرین با تنظیمات مختلف: S1. تجزیه و تحلیل گفتار و ساخت ویژگی: S2. تراز قاب از بدنه موازی: S3. آموزش تابع تبدیل. برداری برای همه شرایط فعال است.

گفته ها (برای تراز قاب) می توانند به طور مستقل و همزمان در هسته های جداگانه پردازش شوند. برای بررسی بیشتر توزیع بار کار، ما زمان اجرای هر مرحله در طول تمرین را با ۲۰ نمونه تمرین موازی، که هر کدام به طور متوسط حدود ۴,۵ ثانیه طول کشید، اندازه گیری کردیم. نتایج در شکل ۱۰ (a) ترسیم شده است.

به وضوح نشان داده شده است که مرحله اول، یعنی تجزیه و تحلیل گفتار و ساخت ویژگی، بیشتر زمان می برد (حدود ۸۴٪). در مقابل، مرحله دوم آموزش، یعنی هم تراز فریم، تنها نسبت ناچیزی از کل زمان مجموعه آموزشی فعلی را می گیرد. در مورد مرحله آخر، آموزش تابع تبدیل، اگرچه زمان محاسباتی غیر قابل اغماض را می طلبد، موازی کردن این مرحله ساده نیست، زیرا به کل داده ها نیاز دارد. علاوه بر این، از آنجایی که آخرین مرحله در مقایسه با دو مرحله اول زمان بسیار کمتری می برد، ما در حال حاضر تنها دو مرحله اول را از طریق موازی سازی چند هسته ای اجرا می کنیم. با توزیع حجم کار دو مرحله اول به چهار هسته با استفاده از چهار رشته، زمان اجرای این دو مرحله به میزان زیادی به حدود ۳۵,۲٪ و کل زمان اجرا به حدود ۴۴,۸٪ از زمان مورد نیاز روی یک هسته کاهش می یابد. که در شکل ۱۰ (ب) نشان داده شده است. در نهایت، می خواهیم تأکید کنیم که اگرچه زمان صرف شده توسط مرحله هم تراز فریم (S2) در شکل ۱۰ (جزئی به نظر می رسد، اما همچنان لازم است این مرحله را موازی کنیم، زیرا پیچیدگی زمانی نظری آن به دلیل استفاده از تاب خوردگی زمانی پویا [۵۳]، در حالی که دو مرحله دیگر فقط پیچیدگی زمانی خطی را به صورت تجربی نشان می دهند.

(3) اثر تعداد هسته ها برای تجزیه و تحلیل بیشتر اثر افزایش سرعت ناشی از موازی سازی چند هسته ای، زمان آموزش مجموعه آموزشی متشکل از ۲۰ گفته موازی با مدت متوسط حدود ۴,۵ ثانیه و همچنین زمان تبدیل یک گفتار ورودی منبع ۵ ثانیه را اندازه گیری کردیم. از آنجایی که در کل ۴ هسته در CPU

رایانه شخصی ما وجود دارد، برنامه به گونه ای پیکربندی شده است که از تعداد هسته های متفاوتی از ۱ تا ۴ برای موازی سازی استفاده کند. زمان اجرا با توجه به تعداد هسته های مجاز گزارش شده است



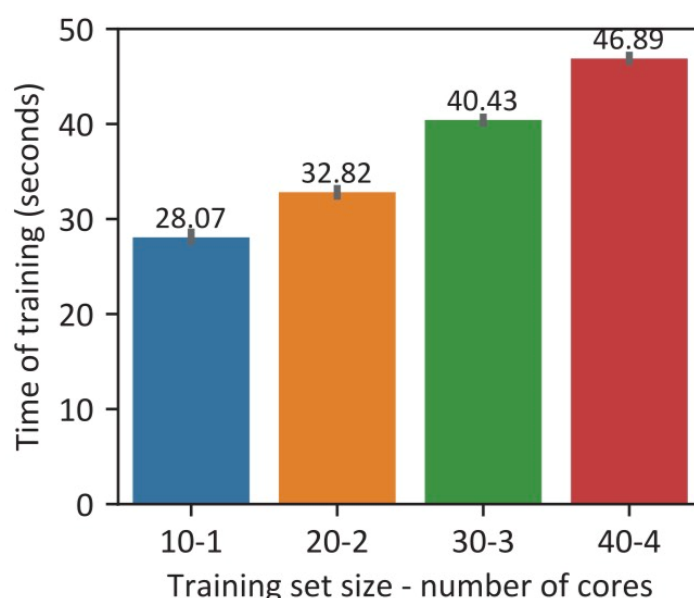
**تصویر ۱۱:** زمان اجرای مرحله آموزش (بالا) و مرحله تبدیل (پایین) زمانی که تعداد متفاوتی از هسته های CPU برای موازی سازی چند هسته ای استفاده می شود (فاصله اطمینان ۹۵). (در کل چهار هسته در CPU در دست بررسی وجود دارد. برداری فعال است.

در شکل ۱۱. از آنجایی که فقط دو مرحله اول مرحله آموزش از موازی سازی چند هسته ای بهره می برند، یعنی حدود ۸۴ میکرومتر از برنامه را می توان موازی کرد (شکل ۱۰ (a))، اثر شتاب نشان داده شده است. شکل ۱۱ در واقع با قانون Amdahl سازگار است، که اغلب در محاسبات موازی برای پیش بینی سرعت تئوری در هنگام استفاده از پردازنده های متعدد استفاده می شود [۵۴]. به طور خاص، با چهار هسته موجود، سرعت مرحله آموزش حدود ۲٫۳ برابر است، در حالی که قانون Amdahl سرعت تئوری را حداکثر ۲٫۷

برابر پیش بینی می کند. به طور کلی، با توجه به رابطه بین زمان اجرا و تعداد هسته های نشان داده شده در شکل ۱۱، این سیستم تبدیل صدا با افزایش سرعت زیر خطی زمانی که موازی سازی چند هسته ای اعمال می شود مشخص می شود.

ما بیشتر مقیاس پذیری موازی سازی چند هسته ای را در سیستم خود با افزایش اندازه مجموعه آموزشی متناسب با تعداد هسته های استفاده شده بررسی کردیم. نتایج در شکل ۱۲ نشان داده شده است. بدیهی است که میانگین زمان مورد نیاز برای ۱۰ نمونه آموزشی در هر هسته با بزرگ شدن اندازه کل مشکل افزایش می یابد. این به دلیل این واقعیت است که همه بخش های چارچوب الگوریتم موازی نیستند (شکل ۱۰). با این وجود، استفاده از موازی سازی چند هسته ای در سیستم ما هنوز مقیاس پذیر است، زیرا زمان اجرا می تواند برای اندازه مشکل ثابت تا حد زیادی کاهش یابد که هسته های بیشتری در دسترس هستند (شکل ۱۱). از منظر عملی، باید توجه داشت که اگرچه اصولاً هنگام اجرای موازی سازی مبتنی بر multithreading می توانیم از تعداد رشته هایی که دوست داریم استفاده کنیم، زمانی که تعداد Thread ها نسبت به هسته ها در یک CPU منفرد بیشتر باشد، سرعت افزایش نمی تواند بیشتر شود. به عنوان مثال، در مورد ما، مرحله آموزش سیستم ما به CPU فشرده است، و بنابراین حداکثر چهار رشته را می توان به طور همزمان روی یک CPU 4 هسته ای اجرا کرد.

تصویر ۱۲: زمان آموزش اندازه های مختلف مجموعه آموزشی و تعداد هسته ها (فاصله اطمینان ۹۵٪).



طور متوسط، هر هسته مربوط به یک مجموعه آموزشی از ۱۰ بیان است. برچسب  $n-k$  به معنای  $n$  گفته در مجموعه آموزشی و  $k$  هسته است. برداری فعال است.

## نمایش برنامه

برنامه موبایلی که ما برای iPhone 7 توسعه داده ایم، Voichap، در مجموع دارای شش پنجره است که عکس های صفحه نمایش آنها در شکل ۱۳ نشان داده شده است که با مازول های مختلف در شکل ۷ مطابقت دارد. مهمترین اقدامات در شکل ها انجام شده است. ۱۳ (d) و ۱۳ (f)، جایی که آموزش و تبدیل

انجام می شود. همانطور که در شکل ۱۳ (d) نشان داده شده است، ما تعدادی رونوشت را برای خواندن توسط گوینده منبع و گوینده هدف برای ساختن یک پیکره موازی برای اهداف آموزشی ارائه کرده ایم، جایی که می توانیم اطلاعات گوینده هدف، از جمله او را نیز ویرایش کنیم. نام او و توضیحات کوتاه به طور خاص، در زمینه Voichap، سخنران منبع در واقع به کاربر اشاره دارد. از این رو، کاربر فقط باید این جملات را یک بار بخواند و صداهای ذخیره شده را می توان دوباره برای مطابقت با همه بلندگوهای هدف برای آموزش در آینده استفاده کرد. پس از آن، کاربر می تواند از گوینده هدف، به عنوان مثال، دوست خود، بخواهد که این جملات را دوباره روی تلفن بخواند و گفته ها توسط این برنامه ضبط می شود تا مجموعه آموزشی ایجاد شود.

همانطور که قبلاً در شکل ۷ نشان داده شد، اطلاعات بلندگو در یک پایگاه داده سبک ذخیره می شود و فایل های صوتی ضبط شده در یک فهرست خاص برای هر بلندگو در سیستم فایل محلی iOS سازماندهی می شوند. پس از اتمام ضبط، دکمه Train در شکل ۱۳ (d) را می توان برای شروع فرآیند آموزشی لمس کرد و در پایان، یک فایل مدل تبدیل مربوط به جفت منبع-هدف داده شده ایجاد می شود. توجه داشته باشید که فایل مدل حاوی تمام اطلاعات لازم برای بازیابی سریع تابع تبدیل  $f(\cdot)$  در صورت نیاز است (شکل ۱). اکنون، زمان بازی با تبدیل صدای بلادرنگ است که توسط صفحه Speak here در شکل ۱۳ (f) مشخص شده است. فقط کافی است دکمه میکروفون را فشار داده و نگه دارید و هر چیزی که دوست دارید بگویید. پس از رها شدن دکمه میکروفون، مرحله تبدیل به صورت خودکار و بلافاصله شروع می شود.

پس از مدت کوتاهی، می توانید آنچه را که اخیراً گفته اید تکرار می شود، اما با صدای گوینده هدف بشنوید. به عنوان یک یادداشت جانبی، برنامه Voichap ما همچنین می تواند با فایل های صوتی به طور مستقیم به غیر از ضبط جملات بداهه توسط خودمان تغذیه شود. این پشتیبانی می تواند برنامه را بسیار جالب تر کند. همانطور که در بخش مقدمه ذکر شد، در حال حاضر، برخی از برنامه های پولی در گوگل پلی یا اپل استور در رابطه با تغییر صدا در دسترس هستند که می توانند صدای افراد مشهور را تقلید کنند، اگرچه عملکرد آنها به طور کلی ضعیف است. با این حال، با برنامه Voichap ما، می توانید صدای خود را به صدای هر فرد مشهوری تبدیل کنید، زمانی که بتوانید برخی از گفته های آموزشی از او دریافت کنید. البته این غیرعملی است که واقعاً از یک سلبیتهی بخواهیم متن مشخص شده را روی تلفن ما بخواند. با این حال، به راحتی می توانیم سخنرانی های عمومی آنها را از وبسایت هایی مانند یوتیوب دانلود کرده و سپس فایل های صوتی را استخراج کرده و به عنوان مواد آموزشی سخنران مورد نظر به جملات تقسیم کنیم. کار باقی مانده این است که همان جملات را بیان کنیم و گفته هایمان را به عنوان داده های آموزشی گوینده منبع ضبط کنیم. این فرآیند در شکل ۱۴ نشان داده شده است، جایی که ما رئیس جمهور ترامپ را به عنوان مثال در نظر می گیریم. به طور خلاصه، ابتدا یک هدف جدید ایجاد می شود و فایل های صوتی بارگذاری می شوند که ما آنها را از YouTube برداشتیم. پس از آموزش مدل، می توانیم کلماتی مانند «آمریکا را دوباره بزرگ کن» به زبان بیاوریم که گویی رئیس جمهور ترامپ هستیم، در حالی که سیستم تبدیل صدا، فردیت صدا را تغییر می دهد و آن را طوری می سازد که گویی واقعاً توسط رئیس جمهور ترامپ گفته شده است.

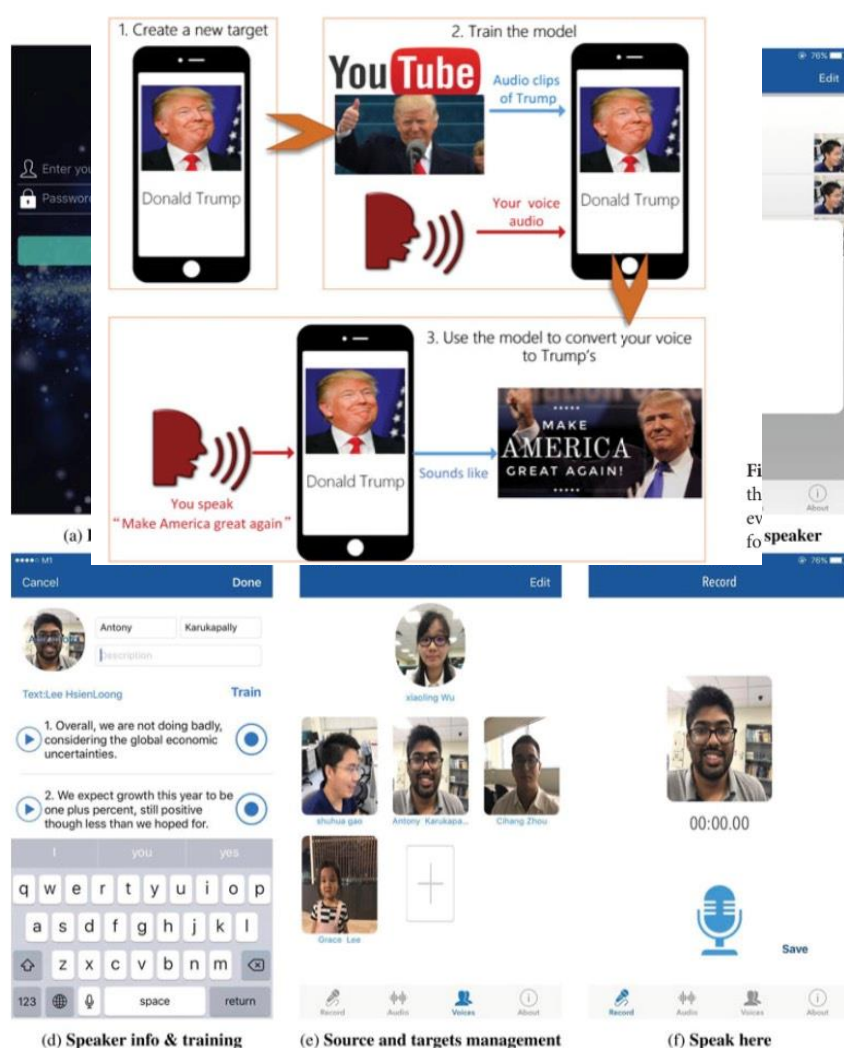
## آزمایشی ارزیابی نرم افزار سیستم تبدیل صدا

اثربخشی و عملکرد بلادرنگ سیستم تبدیل صدای تلفن همراه ما، Voichap، با تبدیل صدا در میان سخنرانان زن و مرد مورد ارزیابی قرار گرفت. برخلاف معیار سنتی تبدیل صدا، که همیشه به دنبال تشابه تبدیل بهتر یا کیفیت گفتار بدون توجه به هزینه محاسباتی است، به عنوان یک برنامه کاربردی موبایل، Voichap تلاش می‌کند تا در یک بازه زمانی قابل قبول به بهترین عملکرد اما نه لزوماً بهترین عملکرد برسد. هدف نهایی ما توسعه یک سیستم تبدیل صدای بلادرنگ در تلفن های همراه برای استفاده روزانه است و تجربه کاربر بیشترین اهمیت را دارد. بدیهی است که در این شرایط اگر کاربر مجبور باشد برای دریافت گفتار تبدیل شده نیم ساعت صبر کند، حتی اگر کیفیت تبدیل بسیار بالا باشد، چندان منطقی نیست.

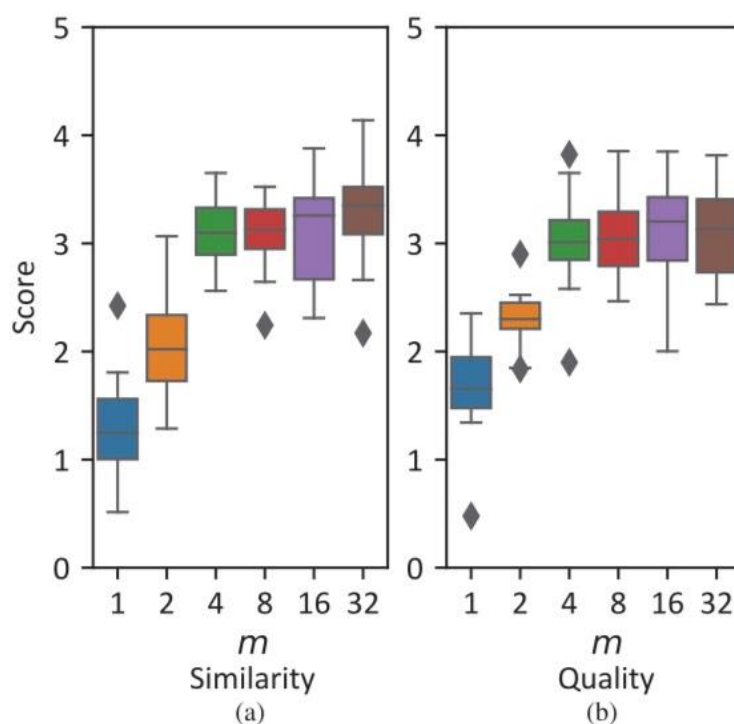
چهار نامزد در محدوده سنی ۲۰ تا ۳۰ سال به عنوان سخنران در آزمون های زیر انتخاب شدند، که شامل یک مرد و یک زن به عنوان منبع و تنظیم یکسان برای هدف بود. بنابراین، چهار جفت گوینده در مجموع بررسی می‌شوند: مرد به مرد (M-M)، مرد به زن (M-F)، زن به مرد (F-M) و زن به زن (F-F).  
(1) تنظیم پارامتر و بازده زمانی

دو پارامتر اصلی سیستم تبدیل صدای ما، تعداد جملات آموزشی و تعداد اجزای گاوسی در مدل GMM، یعنی  $m$  در معادله (۷) است. بر اساس نظرسنجی از چندین کاربر، ما دریافتیم که ساخت یک مجموعه رونویسی از ۲۰ جمله با مدت زمان حدود ۴ یا ۵ ثانیه به منظور آموزش، انتخاب خوبی است. اگرچه به طور کلی یک مجموعه آموزشی بزرگتر می تواند به عملکرد بهتر کمک کند، ممکن است کاربر را از خواندن جملات زیاد خسته کند و همچنین می تواند منجر به زمان طولانی آموزش شود. برای تعیین بهترین گزینه  $m$ ، عمدتاً باید بین زمان اجرا و عملکرد تبدیل تعادل برقرار کنیم. با آزمایش مقادیر مختلف  $m$ ، زمان تمرین را با توجه به هر  $m$  در شکل ۱۵ خلاصه می کنیم.

واضح است که اگر مولفه های گاوسی بیشتری به کار گرفته شوند، مرحله آموزش زمان بیشتری طول خواهد کشید. علاوه بر این، می بینیم که کل زمان آموزش با توجه به تعداد مولفه های گاوسی  $m$  رابطه غیرخطی نشان می دهد. از طرف دیگر، میانگین زمان تمرین برای مجموعه تمرینی شامل ۲۰ گفته بین ۲۰ تا ۵۵ ثانیه برای  $m$  های مختلف است. مهمتر از همه، از طریق یک آزمون گوش دادن ذهنی از گفتار تبدیل شده از نظر تشابه تبدیل شده به هدف و کیفیت گفتار، که در شکل ۱۶ نشان داده شده است، متوجه می شویم که بهبود کلی یک بار  $m > 4$  جزئی است و در نتیجه شنوندگان نظرات قابل مقایسه ای در مورد شباهت و کیفیت گفتارهای تبدیل شده تولید شده توسط مقادیر مختلف  $m$  بالای این آستانه ارائه می دهند. بنابراین، از شکل های ۱۵ و ۱۶، انتخاب  $m = 4$  به عنوان مقدار پارامتر پیش فرض سیستم ما معقول است، که می تواند به کاهش زمان اجرا کمک زیادی کند، اما بدون اینکه کیفیت گفتار تبدیل شده و شباهت را از جنبه انسانی بدتر کند. ادراک در ارزیابی های عینی و ذهنی زیر، همه ما  $\text{setm.4} =$



تصویر ۱۴: یک سناریو معمولی چگونه صدای خود را شبیه ترامپ کنیم؟



تصویر ۱۴: زمان اجرای مرحله تمرین با توجه به تعداد متفاوت مولفه‌های گاوسی (فاصله اطمینان ۹۵) اندازه‌گیری شد. مجموعه آموزشی شامل ۲۰ جمله با طول متوسط حدود ۴,۵ ثانیه است و زمان اجرا با تکرار ۱۰ بار اندازه‌گیری می‌شود. بردارسازی و موازی سازی ۴ هسته ای فعال هستند.

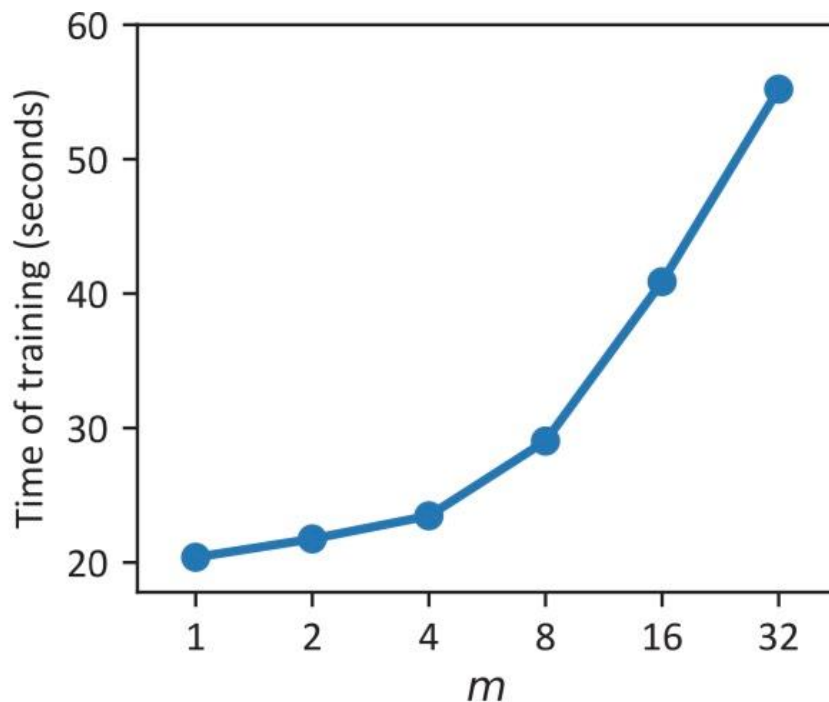
بیایید اکنون به برنامه تلفن همراه Voichap برگردیم. با استفاده از همان مجموعه آموزشی و مجموعه تست، زمان اجرای برنامه با موازی سازی چند هسته ای و چند رشته ای فعال در آیفون ۷ به شرح زیر اندازه‌گیری شد:

- مرحله آموزش یک پیکره موازی شامل ۲۰ جمله با طول متوسط حدود ۴ ثانیه تقریباً ۳۶,۵ ثانیه طول می‌کشد.

- مرحله تبدیل برای ۴,۶ جمله از بلندگوی منبع حدود ۰,۸۵ ثانیه طول می‌کشد.



این نتیجه انتظارات ما را بر اساس شکل‌های ۹ و ۱۵ برآورده می‌کند. توجه داشته باشید که در مقایسه با رایانه شخصی، توانایی محاسباتی دستگاه آیفون نسبتاً ضعیف است، که به طور اجتناب‌ناپذیری منجر به طولانی‌تر شدن زمان اجرا نسبت به زمان اندازه‌گیری شده با MATLAB در رایانه شخصی ویندوز ۱۰



می‌شود\*.

تصویر ۱۶: ارزیابی ذهنی سیستم تبدیل صدا با توجه به تعداد اجزای گاوسی  $m$  در GMM. از ۱۰ شنونده خواسته شد تا شباهت و کیفیت گفتار را ارزیابی کنند. در اینجا میانگین امتیاز چهار جهت تبدیل منبع-هدف ممکن را نشان می‌دهد.

**Table 1.** The MCD of the unconverted source, the traditional GMM and the weighted frequency warping (WFW) method

	No conversion	GMM	WFW
MCD(dB)	7.83	5.85	5.97

(2) ارزیابی عینی تبدیل صدا معمولاً دو نوع ارزیابی را می‌توان برای امتیاز دادن به عملکرد یک سیستم تبدیل صدا استفاده کرد: ارزیابی عینی و ارزیابی ذهنی. برای ارزیابی عینی، پرکاربردترین معیار در ادبیات، اعوجاج Mel-cepstral (MCD) بین گفتار تبدیل شده و گفتار هدف اصلی است [۹، ۱۵، ۱۸]. MCD توسط محاسبه می‌شود

$$\text{MCD}[\text{dB}] = \frac{10}{\log 10} \sqrt{2 \sum_{i=1}^{24} (c_i - \hat{c}_i)^2}, \quad (18)$$

۴- زمان بندی/ گانت چارت:

ردیف	نام فعالیت	زمان/ماه	۱	۲	۳	۴	۵	۶	....	۹
۱	جمع آوری اطلاعات									
۲	بررسی پیشینه									
۳										
۴										
۵										
۶										
۷										
۸										
۹										
۱۰										

نکته: پس از تصویب شورای پژوهشی دانشکده حداقل زمان قابل قبول برای پیش بینی مراحل مطالعاتی و اجرایی پایان نامه کارشناسی ارشد ۶ ماه می باشد.

۵- نظریه شورای گروه تخصصی:

طرح تحقیق پایان نامه خانم / آقای: .....

دانشجوی مقطع کارشناسی ارشد رشته ..... در شورای تخصصی گروه مورخ

..... مطرح شد. پس از بحث و تبادل نظر مورد تصویب اکثریت اعضاء قرار گرفت □ نگرفت □

ردیف	نام و نام خانوادگی	تخصص	نوع رای	امضاء
۱				
۲				
۳				
۴				
۵				

تاریخ:

امضاء:

مدیر گروه:

بسمه تعالی



واحد تهران جنوب حفظ و دفاع از حقوق مادی و معنوی تولیدات علمی دانشگاه آزاد اسلامی و ارائه نتایج آنها  
مربوط با دانشجویان کارشناسی ارشد

عنوان پایان نامه: توسعه یک سیستم تبدیل صوتی کارآمد محاسباتی در تلفن های همراه

نام: پانیز نام خانوادگی: اسلامی تمیجانی شماره دانشجویی: ۴۰۱۱۴۱۴۰۱۱۰۲۶

دانشکده: فنی و مهندسی رشته تحصیلی: مهندسی پزشکی گرایش: بیوالکتریک

سال اخذ پایان نامه: ۱۴۰۱ نیمسال تحصیلی: اول

تلفن: تلفن همراه: پست الکترونیک: paniz.eslami99@gmail.com

تعهدات دانشجو:

- 1- محتوای پایان نامه کارشناسی ارشد، از آن دیگران نیست (دست اول است)، براساس اصول علمی تهیه شده است و با نام دانشگاه آزاد اسلامی- واحد تهران جنوب ارائه خواهند شد.<sup>۱</sup>
- 2- به منظور رجوع مناسب و روشن به آثار دیگران، منابع و مأخذ مربوط به نقل قول ها، جدول ها و نمودارها و یا نتایج تحقیقات دیگران در پایان نامه دقیقاً ذکر خواهد شد؛ همچنین هیچ گونه استفاده ای از آثار دیگران بدون ذکر منبع اصلی و به گونه ای که قابل تشخیص و تفکیک از متن اصلی نباشد، به عمل نخواهد آمد.
- 3- بدون ذکر نام دانشگاه آزاد اسلامی- واحد تهران جنوب و در نظر گرفتن حقوق این دانشگاه، در مورد ارائه و انتشار نتایج حاصل از پایان نامه به شکل مقاله، کتاب، اختراع، اکتشاف و ... (در قالب مطالب چاپی یا غیر چاپی) در هر مرحله (قبل و بعد از دفاع از پایان نامه)، اقدامی صورت نخواهد گرفت. بدیهی است که ارسال هر مقاله مستخرج از پایان نامه باید با هماهنگی با استاد راهنما باشد.
- 4- برای جلوگیری از درج مقاله در نشریات بی اعتبار، قبل از چاپ مقاله، اعتبار نشریه از فهرست نشریات بی اعتبار در سایت معاونت پژوهشی و فناوری دانشگاه آزاد اسلامی به نشانی <http://sp.rvp.iau.ir> بررسی خواهد شد.
- 5- در صورت هرگونه مغایرت و تخلف از موارد اشاره شده در بندهای ۱ تا ۳ این تعهدنامه، دانشگاه آزاد اسلامی- واحد تهران جنوب مجاز است از ادامه تحصیل و هرگونه فعالیت آموزشی و امکان دفاع از پایان نامه دانشجو در هر مرحله از تحصیل جلوگیری کند. همچنین خسارات مادی و معنوی وارده به دانشگاه آزاد اسلامی و افراد ذی نفع پرداخت خواهد شد.

نام و نام خانوادگی دانشجو:

امضاء:

تاریخ:

مقالاتی تحت بررسی قرار خواهند گرفت که طبق بخشنامه های سازمان مرکزی باشند.

1- بخشنامه شماره ۷۲/۳۴۵۱۹ مورخ ۹۲/۲/۱۲ باشد. مفاد بخشنامه .... "در صورتی که نام فرد دیگری به غیر از استاد راهنما، مشاور و دانشجو در تیم نویسندگان مقاله مستخرج از پایان نامه و رساله ها قید گردد؛ به مقاله مذکور در مقطع کارشناسی ارشد و دکترای حرفه ای نمره ای اختصاص نمی یابد...."

2- بخشنامه شماره ۷۲/۲۹۹۹۲۰ مورخ ۹۲/۹/۹ باشد. مفاد بخشنامه: ".... در مقاله های مستخرج، نویسنده اول دانشجو و به نام واحد تحصیل دانشجو و استاد راهنما عهده دار مکاتبات است...."

3- بخشنامه شماره ۷۰/۸۱۲۴۸ مورخ ۹۳/۹/۱ باشد. مفاد بخشنامه "نحوه آدرس دهی

مقاله های انگلیسی: Department of ..., South Tehran Branch, Islamic Azad University, Tehran, Iran

\*توجه: تشخیص نشریات بی اعتبار: دو مورد اصلی در تشخیص نشریات بی اعتبار عبارتند از: ۱- تقاضای اخذ وجه توسط ناشر در زمان ارسال یا پذیرش مقاله و ۲- آدرس

الکترونیکی نشریات بی اعتبار (که اغلب پست های الکترونیکی رایگان نظیر سایت Yahoo و غیره است). همچنین کنترل نشریه در سایت <http://sp.rvp.iau.ir>



واحد تهران جنوب نامه:

تبدیل صوتی کارآمد محاسباتی در تلفن های همراه

حفظ و دفاع از حقوق مادی و معنوی تولیدات علمی دانشگاه آزاد اسلامی و ارائه نتایج آنها

الف) استاد راهنما:

اینجانب استاد راهنمای آقای/ خانم دانشجوی مقطع کارشناسی ارشد دانشگاه آزاد اسلامی- واحد تهران جنوب، از مفاد بخشنامه «حفظ و دفاع از حقوق مادی و معنوی تولیدات علمی دانشگاه آزاد اسلامی و ارائه نتایج آنها»، آگاهی کامل داشته و خود را ملزم به رعایت آن می دانم.

پست الکترونیک:

تلفن:

امضاء:

تاریخ:

ب) استاد مشاور: (در صورت لزوم)

اینجانب استاد مشاور آقای/ خانم دانشجوی مقطع کارشناسی ارشد دانشگاه آزاد اسلامی- واحد تهران جنوب، از مفاد بخشنامه «حفظ و دفاع از حقوق مادی و معنوی تولیدات علمی دانشگاه آزاد اسلامی و ارائه نتایج آنها»، آگاهی کامل داشته و خود را ملزم به رعایت آن می دانم.

پست الکترونیک:

تلفن:

امضاء:

تاریخ:

باسمه تعالی

نامه کارشناسی ارشد



فر

\*(لطفاً در این قسمت چیزی ننویسید.)

واحد تهران جنوب

محل در

مشخصات دانشجو:	
نام و نام خانوادگی دانشجو: .....	شماره دانشجویی: .....
مجتمع/دانشکده: .....	
رشته تحصیلی: .....	تعداد واحد پایان نامه: .....
پایان نامه: اول ..... / دوم .....	
امضاء کارشناس آموزش مجتمع/ دانشکده: .....	
امضاء رئیس اداره آموزشی مجتمع/ دانشکده: .....	
عنوان پایان نامه:	
نام و نام خانوادگی استاد راهنما:	
رشته تحصیلی:	مرتبه علمی:
نوع همکاری: تمام وقت <input type="checkbox"/> نیمه وقت <input type="checkbox"/>	عضو هیات علمی مدعو از سایر واحدهای دانشگاه آزاد اسلامی <input type="checkbox"/>
عضو هیات علمی مدعو از دانشگاه دولتی <input type="checkbox"/>	عضو غیر هیات علمی <input type="checkbox"/>
امضاء استاد:	
نام و نام خانوادگی استاد مشاور:	
رشته تحصیلی:	مرتبه علمی:
نوع همکاری: تمام وقت <input type="checkbox"/> نیمه وقت <input type="checkbox"/>	عضو هیات علمی مدعو از سایر واحدهای دانشگاه آزاد اسلامی <input type="checkbox"/>
عضو هیات علمی مدعو از دانشگاه دولتی <input type="checkbox"/>	عضو غیر هیات علمی <input type="checkbox"/>
امضاء استاد:	
نام و نام خانوادگی مدیر گروه آموزشی – پژوهشی .....	
تاریخ و امضاء	
تاریخ تصویب پایان نامه در شورای پژوهشی مجتمع/دانشکده: ..... شماره جلسه: .....	

نکته ۱: تمام اطلاعات این فرم صحیح و کامل تایید شود و به تایید اساتید مربوطه رسانده شود.

نکته ۲: ارسال تصویر کارت ملی (پشت و رو)، آخرین حکم هیئت علمی، رزومه علمی، آخرین مدرک تحصیلی برای کلیه استادان راهنما و مشاور مدعو (عضو هیئت علمی سایر واحدهای دانشگاه آزاد اسلامی و یا وزارتین) برای یک بار الزامی است.

نکته ۳: مسئولین مربوطه می بایست اصل این فرم را به همراه صورتجلسات پروپوزال های تصویب شده در شورای پژوهشی مجتمع/ دانشکده و فرم شماره ۱ فایل (Excel) را بطور همزمان به حوزه معاونت پژوهش و فناوری واحد ارسال نمایند.



فرم تصویب (پروپوزال) مربوط به دانشجو ..... به شماره دانشجویی .....  
 رشته ..... در تاریخ ..... در شورای  
 پژوهشی مجتمع/دانشکده مطرح و تصویب گردید.

این طرح در تاریخ ..... در شورای پژوهشی مجتمع/دانشکده مطرح گردید ولی به علل زیر مورد موافقت قرار نگرفت.

علل عدم تصویب طرح تحقیق پایان نامه (پروپوزال):