

صلى الله عليه وسلم



دانشکده برق و رباتیک

رشته‌ی مهندسی برق گرایش مخابرات سیستم

پایان نامه کارشناسی ارشد

تشخیص مرز کلمات درسیگنال گفتار با استفاده از روش زمان – فرکانس

نگارنده:

فاطمه تاجیک ایجدان

استاد راهنما:

دکتر حسین مروی

شهریور ۱۴۰۰

تقدیم به

پدر و مادر عزیز و مهربانم که در سختی‌ها و دشواری‌های زندگی همواره یوری

دلسوز و فداکار و پشتیبانی محکم و مطمئن برایم بوده‌اند.

سپاس بی کران پروردگار یکتا را که هستی مان بخشید و به طریق علم و دانش رهنمونمان
شد و به همنشینی رهروان علم و دانش مفتخرمان نمود و خوشه چینی از علم و معرفت را
روزیمان ساخت.

بدینوسیله از زحمات و تلاش بی دریغ خانواده عزیز و استاد محترم جناب آقای دکتر
حسین مروی که در تهیه این مجموعه با این جانب همکاری داشته اند، تشکر و مراتب سپاس
قلبی خود را اعلام نموده و موفقیت ایشان را از خداوند متعال خواهانم.

تعهدنامه

اینجانب **فاطمه تاجیک ایجدان** دانشجوی دوره کارشناسی ارشد رشته مخابرات سیستم دانشکده برق و ریاتیک، دانشگاه صنعتی شاهرود نویسنده پایان نامه با موضوع **تشخیص مرز کلمات در سیگنال گفتار با استفاده از روش زمان - فرکانس** تحت راهنمایی دکتر **حسین مروی** متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورد استفاده استناد شده است.
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می‌باشد و مقالات مستخرج با نام **دانشگاه صنعتی شاهرود** و یا **Shahrood University of Technology** به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می‌گردد.
- در کلیه مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت‌های آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ:

امضاء دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات، مستخرج، کتاب، برنامه‌های رایانه ای، نرم‌افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی‌باشد.

چکیده

تشخیص مرز کلمه یعنی شناسایی آغاز و پایان کلمه در یک سیگنال گفتاری که نقش مهمی برای بسیاری از کاربردها مانند تشخیص کلمه کلیدی، تشخیص گفتار، سیستم تبدیل متن به گفتار و ... دارد. در سیستم تشخیص گفتار بیشتر خطاها ناشی از تشخیص اشتباه مرزهای کلمات است. تشخیص صحیح مرز کلمه، خطاها را کاهش می‌دهد و عملکرد سیستم‌ها را بهبود می‌بخشد. الگوریتم‌های زیادی برای تشخیص مرز کلمات استفاده می‌شود که برخی از آن‌ها به دلیل نویز محیط و یا عدم مکث خوب بین کلمات عملکرد چندان خوبی برای تشخیص صحیح مرز کلمات ندارند. در صورتی که برخی از روش‌ها با هم ترکیب شوند بهتر می‌توانند مرزهای کلمات را تشخیص دهند.

در این تحقیق هدف تشخیص مرزهای کلمه با استفاده از روش‌های زمان-فرکانس از جمله موجک و ویگنر-ویل است. روش‌های زمان-فرکانس سیگنال را در هر دو حوزه زمان و فرکانس بررسی می‌کنند. چون سیگنال گفتار یک سیگنال ناپایستا است و مشخصات آن با زمان تغییر می‌کند بهترین روش برای تحلیل این سیگنال و استخراج ویژگی از آن روش‌های زمان-فرکانس می‌باشد. روش‌های مختلف استخراج ویژگی مانند MFCC، LPC، PLP، تحلیل موجک و روش ویگنر-ویل جهت تشخیص مرز کلمات مورد بررسی قرار می‌گیرند. از دونه داده برای بررسی نتایج استفاده شده است. در پایگاه داده TIMIT زمان بین کلمات خیلی کم است و به طور طبیعی و تقریباً سریع کلمات تلفظ می‌شوند، در حالی که در CADLab کلمات به صورت شمرده بیان می‌شوند. روش ویگنر-ویل قادر به تشخیص مرزهای کلمات با ۸۳/۵٪ برای داده اول، و ۶۸٪ برای داده دوم بهترین روش برای شناسایی مرزهای کلمات می‌باشد. آنتروپی نیز با ۲۵٪ تشخیص برای داده اول و نرخ عبور از صفر با ۳۰٪ تشخیص برای داده دوم کمترین میزان تشخیص صحیح مرزهای کلمات را در میان روش‌های انجام شده دارد.

واژگان کلیدی: تشخیص مرز کلمه، استخراج ویژگی، زمان-فرکانس، MFCC، PLP، LPC

لیست مقالات مستخرج از پایان نامه:

...

فهرست

- ۱- فصل اول: ۱
- ۱-۱- مقدمه ۲
- ۲-۱- استخراج ویژگی ۲
- ۳-۱- تشخیص مرز کلمه ۴
- ۴-۱- تقسیم‌بندی گفتار به بخش‌های واکدار و بی‌واک ۵
- ۵-۱- ضرورت انجام پژوهش ۶
- ۶-۱- ساختار پایان‌نامه ۷
- ۲- فصل دوم: ۹
- ۱-۲- مقدمه ۱۰
- ۲-۲- انرژی زمان-کوتاه ۱۰
- ۱-۲-۲- الگوریتم کلی از انرژی زمان-کوتاه ۱۱
- ۲-۲-۲- بررسی کارهای گذشته ۱۱
- ۳-۲- نرخ عبور از صفر ۱۲
- ۱-۳-۲- نشانه تشخیص مرزها با استفاده از نرخ عبور از صفر ۱۳
- ۲-۳-۲- بررسی کارهای گذشته ۱۳
- ۴-۲- گام ۱۵
- ۱-۴-۲- الگوریتم تشخیص مرز کلمه با استفاده از روش گام ۱۶
- ۲-۴-۲- انجام کارهای گذشته ۱۶
- ۵-۲- آنتروپی ۱۷
- ۱-۵-۲- نشانه‌های تشخیص مرز کلمه با استفاده از روش آنتروپی ۱۸
- ۲-۵-۲- بررسی کارهای گذشته ۱۸
- ۳- فصل سوم: ۲۳

۲۴	۱-۳- مقدمه
۲۴	۲-۳- مقیاس مل
۲۵	۱-۲-۳- استخراج ویژگی ضرایب کپسترال در مقیاس مل
۲۸	۳-۳- تحلیل پیش‌بینی خطی
۲۹	۴-۳- مقیاس بارک
۳۰	۱-۴-۳- تحلیل ادراکی پیشگویی خطی
۳۳	۵-۳- تبدیل موجک
۳۶	۶-۳- توزیع ویگنر
۳۹	۷-۳- ویژگی‌های توزیع ویگنر-ویل
۴۳	۴- فصل چهارم:
۴۴	۱-۴- مقدمه
۴۴	۲-۴- پایگاه داده
۴۵	۳-۴- پیش‌پردازش
۴۶	۴-۴- نتایج شبیه‌سازی روش‌های متداول
۴۶	۱-۴-۴- نتایج شبیه‌سازی روش انرژی زمان-کوتاه
۴۹	۲-۴-۴- نتایج شبیه‌سازی روش نرخ عبور از صفر
۵۱	۳-۴-۴- نتیجه شبیه‌سازی با استفاده از روش گام
۵۳	۴-۴-۴- نتیجه شبیه‌سازی روش آنتروپی
۵۵	۵-۴- بلوک دیاگرام کلی روش‌های پیشنهادی
۵۶	۱-۵-۴- نتیجه شبیه‌سازی روش MFCC
۵۷	۲-۵-۴- نتیجه شبیه‌سازی روش LPC
۵۹	۳-۵-۴- نتیجه شبیه‌سازی روش PLP
۶۰	۴-۵-۴- تئوری کار روش موجک
۶۳	۵-۵-۴- تئوری کار روش ویگنر-ویل

۶-۴- نتیجه‌گیری ۶۵

۵- فصل پنجم: ۶۹

۵-۱- جمع‌بندی ۷۰

۵-۲- پیشنهادات برای ادامه کار ۷۱

۶- منابع ۷۳

فهرست شکل‌ها

- شکل ۱-۲: بلوک دیاگرام تشخیص گام با استفاده از AMDF و ACF [۱۴] ۱۶
- شکل ۲-۲: الگوریتم تشخیص گفتار با استفاده از روش آنروپی و انرژی زیرباند [۱۹] ۲۰
- شکل ۳-۲: تشخیص مرزها تحت نویز سفید [۱۹] ۲۱
- شکل ۴-۲: تشخیص مرزها تحت نویز همهمه [۱۹] ۲۱
- شکل ۱-۳: تغییرات مقیاس مل برحسب رابطه (۳-۱) [۲۵] ۲۵
- شکل ۲-۳: مراحل استخراج ویژگی MFCC [۲۲] ۲۶
- شکل ۳-۳: پاسخ فرکانسی فیلتربانک مل [۲۵] ۲۷
- شکل ۴-۳: تغییرات مقیاس بارک نسبت به رابطه (۳-۸) [۲۸] ۳۰
- شکل ۵-۳: مراحل استخراج ویژگی PLP ۳۱
- شکل ۶-۳: فیلتربانک بارک [۳۰] ۳۲
- شکل ۷-۳: تبدیل موجک گسسته با استفاده از فیلتر بانک [۳۴] ۳۵
- شکل ۸-۳: نمودار مراحل تجزیه سیگنال به سه سطح [۳۴] ۳۵
- شکل ۹-۳: مراحل تجزیه سیگنال به سه سطح [۳۴] ۳۶
- شکل ۱۰-۳: سیگنال چیرپ [۴۰] ۳۷
- شکل ۱۱-۳: تبدیل ویگنر-ویل سیگنال چیرپ [۴۰] ۳۷
- شکل ۱۲-۳: وجود عبارات متقاطع در توزیع ویگنر [۳۷] ۳۸
- شکل ۱-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه ۴۷
- شکل ۲-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۴۰ نمونه ۴۷
- شکل ۳-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و

- شيفت زمانی ۸۰ نمونه ۴۸
- شکل ۴-۴: تشخيص مرز کلمه با استفاده از روش انرژی برای داده دوم با طول پنجره ۳۲۰ نمونه و شيفت زمانی ۲۰ نمونه ۴۸
- شکل ۴-۵: تشخيص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شيفت زمانی ۲۰ نمونه ۵۰
- شکل ۴-۶: تشخيص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شيفت زمانی ۴۰ نمونه ۵۰
- شکل ۴-۷: تشخيص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شيفت زمانی ۸۰ نمونه ۵۰
- شکل ۴-۸: تشخيص مرز کلمه با استفاده از نرخ عبور از صفر برای داده دوم با طول پنجره ۳۲۰ نمونه و شيفت زمانی ۲۰ نمونه ۵۱
- شکل ۴-۹: تشخيص مرز کلمه با استفاده از روش گام برای داده اول ۵۲
- شکل ۴-۱۰: تشخيص مرز کلمه با استفاده از روش گام برای داده دوم ۵۳
- شکل ۴-۱۱: تشخيص مرز کلمه با استفاده از روش آنتروپی برای داده اول ۵۴
- شکل ۴-۱۲: تشخيص مرز کلمه با استفاده از روش آنتروپی برای داده دوم ۵۴
- شکل ۴-۱۳: بلوک دیاگرام پیشنهادی برای تشخيص مرز کلمات ۵۵
- شکل ۴-۱۴: تشخيص مرز کلمه با استفاده از روش MFCC برای داده اول ۵۶
- شکل ۴-۱۵: تشخيص مرز کلمه با استفاده از روش MFCC برای داده دوم ۵۷
- شکل ۴-۱۶: تشخيص مرز کلمه با استفاده از روش LPC برای داده اول ۵۸
- شکل ۴-۱۷: تشخيص مرز کلمه با استفاده از روش LPC برای داده دوم ۵۸
- شکل ۴-۱۸: تشخيص مرز کلمه با استفاده از روش PLP برای داده اول ۶۰
- شکل ۴-۱۹: تشخيص مرز کلمه با استفاده از روش PLP برای داده دوم ۶۰

- شکل ۴-۲۰: تشخیص مرز کلمات با استفاده از روش موجک برای داده اول ۶۲
- شکل ۴-۲۱: تشخیص مرز کلمات با استفاده از روش موجک برای داده دوم ۶۳
- شکل ۴-۲۲: تشخیص مرز کلمات با استفاده از روش ویگنر-ویل برای داده اول ۶۵
- شکل ۴-۲۳: تشخیص مرز کلمات با استفاده از روش ویگنر-ویل برای داده دوم ۶۵
- شکل ۴-۲۴: روش به دست آوردن دقت تشخیص مرز کلمات ۶۶
- شکل ۴-۲۵: نمودارهای دقت برای داده اول و داده دوم ۶۸

فهرست جداول

- جدول ۱-۱: تفاوت بین گفتار بی‌واک و واکنار ۶
- جدول ۱-۲: مقایسه حداقل دقت به دست آمده در سه الگوریتم [۱۳] ۱۴
- جدول ۲-۲: مقایسه حداکثر دقت بدست آمده در سه الگوریتم [۱۳] ۱۴
- جدول ۳-۲: مقایسه خطای سه روش برای تشخیص گام [۱۴] ۱۶
- جدول ۱-۴: تقسیم بندی جملات پایگاه داده TIMIT [۴۱] ۴۴
- جدول ۲-۴: شرح پایگاه داده TIMIT [۴۱] ۴۵
- جدول ۳-۴: نتایج روش‌های تشخیص مرز کلمه ۶۷

فصل اول

کلیات پژوهش

۱-۱- مقدمه

گفتار طبیعی‌ترین و ساده‌ترین روش ارتباط بین انسان‌ها و به عنوان یک وسیله ارتباط جمعی می‌باشد که عامل پیدایش زبان‌های مختلف شده است. گفتار از مهم‌ترین سیگنال‌های پردازشی و ارتباطی است و از یک رشته از صداها تشکیل شده است. صداها نماینده اطلاعاتی هستند که بین انسان‌ها مشترک اند. اغلب پدیده‌های فیزیکی در جهان دارای صدا بوده که تحلیل سیگنال این اصوات ممکن است اطلاعات بسیار مفیدی را ارائه دهند. پردازش گفتار یعنی ارتباط بین انسان و ماشین در اساسی‌ترین شکل طبیعی آن. هدف از پردازش گفتار، تحلیل گفتار برای تشخیص اتوماتیک و استخراج اطلاعات، دریافت برخی از اطلاعات فیزیولوژی گوینده، ارتباط گفتاری بین انسان و ماشین (کامپیوتر) در اساسی‌ترین شکل آن^۱ و پردازش انواع مختلف اصوات جهت کسب یکسری اطلاعات مفید و مورد نظر آن‌ها (از سیگنال صوت) توسط ماشین. از کاربردهای پردازش گفتار می‌توان به شناسایی گفتار و گوینده، تبدیل متن به گفتار، فشرده سازی و بهسازی گفتار و... اشاره کرد.

سیگنال صوتی و شنیداری^۲ احساس ارتعاشات و نوسانات هوا توسط گوش انسان است. اگر این صوت در محدوده فرکانسی ۲۰ هرتز تا ۲۰ کیلوهرتز باشد آکوستیک^۳ نام دارد که با رسیدن به گوش و سپس انتقال به مغز و پردازش روی آن باعث درک مفهوم در ذهن انسان می‌شود. سیگنال گفتار^۴ زیر مجموعه-ای از سیگنال صوتی و شنیداری می‌باشد که توسط انسان ایجاد می‌شود.

۱-۲- استخراج ویژگی^۵

استخراج ویژگی یک مرحله مهم در تشخیص مرز کلمات در یک فایل گفتار است. فرآیندی است که بخش‌های مهم اطلاعات را نگه داشته و بخش‌های اضافی را حذف می‌کند. هدف از استخراج ویژگی

^۱Human machine Interaction (HMI)

^۲ Audio

^۳ acoustic

^۴ speech signal

^۵ Feature Extraction

نمایش هر سیگنال گفتاری با تعداد محدودی ویژگی است. دلیل این امر آن است که تمام اطلاعات موجود در سیگنال صوتی برای پردازش بسیار زیاد است و همه اطلاعات برای کارهای خاص مناسب نیستند. پس تبدیلاتی روی سیگنال گفتار انجام می‌شود تا بتوان ویژگی مورد نظر را استخراج کرد [۱]. برای مثال، ممکن است برخی ویژگی‌ها در محیط‌های نویزی جواب بهتری دهند یا اینکه حجم محاسبات کمتری را طلب کنند. این ویژگی‌ها گاه در کل سیگنال گفتار و گاه در یک پنجره از آن محاسبه می‌شوند، ولی در تحلیل سیگنال گفتار به‌طور معمول ویژگی‌های سیگنال پنجره‌بندی شده تحلیل می‌شود. این ویژگی‌ها گاه یک عددند، نظیر انرژی یا میانگین پنجره سیگنال گفتار و گاه یک بردارند. برای مثال، خود اندازه سیگنال گفتار در حوزه‌ی زمان را می‌توان به صورت برداری یک‌بعدی با طول تعداد نقاط پنجره گفتار در نظرگرفت. از این‌رو، گاه از این ویژگی‌ها با نام ضرایب ویژگی یا بردارهای ویژگی نیز یاد می‌شود. از معروف‌ترین این ویژگی‌ها می‌توان به ضرایب کپسترال فرکانس مل ($MFCC^1$)، تعداد عبور از صفر (ZRC^2)، ضرایب پیشگویی خطی (LPC^3) و تحلیل ادراکی پیشگویی خطی (PLP^4) اشاره کرد، که رایج‌ترین ویژگی که در تشخیص گفتار استفاده می‌شود MFCC است [۲].

الگوریتم‌های استخراج ویژگی به سه روش دسته بندی می‌شوند:

(۱) تحلیل زمانی

(۲) تحلیل طیفی

(۳) تحلیل زمان-فرکانس

¹ Mel Frequency Cepstral Coefficients

² Zero Crossing

³ Linear Prediction Coefficients

⁴ Perceptual linear predictive

۱-۳- تشخیص مرز کلمه^۱

تشخیص مرز کلمه به شناسایی آغاز و پایان هر کلمه در یک عبارت گفتاری دارای نویز گفته می‌شود. این نویز ممکن است نویز پس زمینه یا مصنوعات تولید شده توسط گوینده باشد. نویز پس زمینه مشکل تشخیص مرز کلمات را بطور قابل توجهی پیچیده می‌کند. به عنوان مثال، اغلب ابتدای یا انتهای یک کلمه توسط گوینده مانند صدای دهان، لب زدن و تنفس شدید پنهان می‌شود [۳]. ویژگی‌های مورد نیاز یک آشکارساز مرز کلمه ایده آل عبارتند از: قابل اطمینان بودن، مقاومت در برابر نویز، دقت بالا، سازگاری، سادگی. در میان این ویژگی‌ها، استحکام در برابر نویز بهترین و سختترین ویژگی است. تشخیص مرز کلمه یک مسئله پیچیده است. تشخیص سکوت کافی نیست، زیرا همه کلمات مکث‌های واضحی بین خود ندارند. علاوه بر این، برخی از واج‌ها ویژگی‌هایی شبیه به نویز دارند. مدت زمان تلفظ کلمات تا حدی توسط تعداد و طول واج‌هایی که کلمه را تشکیل می‌دهند و چگونگی تلفظ گوینده تعیین می‌شود [۴]. تشخیص مرز کلمه می‌تواند وابسته به گوینده و مستقل از گوینده طبقه‌بندی شود. سیستم‌های وابسته به گوینده در مقایسه با سیستم‌های مستقل از گوینده به راحتی توسعه می‌یابند، زیرا می‌توان اطلاعاتی در مورد ویژگی‌های گفتار یک گوینده مثل گام بدست آورد، اما سیستم‌های مستقل از گوینده باید براساس بسیاری از گویندگان ایجاد شوند که به دلیل تفاوت در نوای گفتار، گام، سن، جنسیت، گویش و .. تشخیص مرز را به چالش می‌کشند.

تشخیص دقیق مرز کلمه به دو دلیل اهمیت دارد:

(۱) صحت و اعتبار شناسایی گفتار به تشخیص دقیق مرز کلمات بستگی دارد.

(۲) تفکیک کلمات از سکوت، پیچیدگی سیستم را کاهش می‌دهد [۵].

تشخیص مرز کلمه معمولاً با استفاده از اطلاعات صوتی گفتار انجام می‌شود و در بسیاری از کاربردها مانند تشخیص کلمه کلیدی، سیستم تبدیل متن به گفتار و سیستم سنتز و تشخیص گفتار به کار

¹ Word Boundary Detection (WBD)

می‌رود. در سیستم تبدیل متن به گفتار مکث‌های به جا و درست در هنگام تولید گفتار تولید شود تا شنونده معنی را به درستی درک کند [۶]. در هنگام استفاده از اطلاعات مرز کلمات، مهم‌ترین نکته دقیق بودن این اطلاعات می‌باشد. اگر مرزها به درستی شناسایی نشوند، اطلاعات غلط باعث افزایش خطا می‌شود.

تشخیص مرز کلمه یک بخش مهم در سیستم تشخیص خودکار گفتار^۱ است و برای عملکرد کارآمد و دقیق سیستم تشخیص خودکار گفتار ضروری است. دلیل اصلی خطاها در سیستم‌های تشخیص خودکار گفتار، تشخیص نادرست مرزهای شروع و پایان کلمات است. برای الگوریتم‌های تشخیص خودکار گفتار ضروری است که بخش‌های گفتار از غیر گفتاری جدا شوند. فناوری تشخیص گفتار می‌تواند سیگنال گفتار را به کمک ماشین مخصوص به متن یا دستور تبدیل کند و این یک تحقیق بسیار فعال در پردازش سیگنال و هوش مصنوعی است. تشخیص دقیق مرز کلمات در یک سیستم تشخیص گفتار به دلیل کاهش محاسبات در مراحل شناسایی و افزایش دقت در سیستم تشخیص گفتار از اهمیت ویژه‌ای برخوردار است [۷].

۱-۴- تقسیم‌بندی^۲ گفتار به بخش‌های واکدار^۳ و بی‌واک^۴

سیگنال گفتار از سه بخش واکدار، بی‌واک و سکوت تشکیل شده است. دانش ویژگی‌های صوتی به ویژه در بخش‌های واکدار و بی‌واک نقش مهمی را در بسیاری از سیستم‌های تجزیه و تحلیل گفتار ایفا می‌کند. تقسیم بندی گفتار به بخش‌های واکدار، بی‌واک و سکوت^۵ برای پردازش ابتدایی لازم است و در بسیاری از موارد مانند سنتز گفتار، تشخیص گفتار و تقویت گفتار کاربرد دارد [۸]. گفتار واکدار

^۱ Automatic Speech Recognition System (ASRS)

^۲ Segmentation

^۳ voiced

^۴ unvoiced

^۵ silence

تقریباً از یک طیف^۱ ثابت تشکیل شده است. فرکانس های ثابت معمولاً هنگام گفتن مصوت ها^۲ ساخته می شوند. ناحیه صوتی گفتار هنگامی ایجاد می شود که چاکنای^۳ لرزان از طریق مجرای صوتی^۴ طنین انداز شود و سبب بیرون رانده شدن توده‌ی هوا توسط این تارها می‌گردد و پالس‌های متناوب در هوا تولید می‌شود. گفتار واکدار برای فهم گفتار از اهمیت بیشتری برخوردار است. در گفتار بی‌واک تارهای صوتی کاملاً از هم باز هستند و هوا به راحتی از میان آن‌ها عبور می‌کند و پالس های غیر متناوب تولید می‌کند. بنابراین سیگنال های صوتی را می توان به دلیل تناوب آن تشخیص و جدا کرد.

براساس ویژگی‌های گفتار واکدار و بی‌واک ذکر شده در جدول (۱-۱)، الگوریتم‌های زیادی برای تشخیص گفتار آوایی و غیر آوایی وجود دارد. مشکلاتی که در هنگام قطعه بندی با آن‌ها مواجه هستیم:

(۱) همپوشی گفتاری واکدار و بی‌واک با یکدیگر، شناسایی نقاط پایانی دقیق آن‌ها را دشوار می‌سازد.

(۲) گفتار با انرژی کم مانند زمزمه می‌تواند به راحتی با سکوت اشتباه گرفته شود [۹].

جدول ۱-۱: تفاوت بین گفتار بی‌واک و واکدار

گفتار بی‌واک	گفتار واکدار
انرژی پایین در هر فریم	انرژی بالا در هر فریم
نرخ بالا عبور از صفر	نرخ پایین عبور از صفر
شبه به نویز	دارای تناوب

۱-۵- ضرورت انجام پژوهش

کلمه قدیمی‌ترین واحد تشخیص مورد استفاده است. زمانی از کلمه به عنوان واحد تشخیص استفاده

1 spectrum
2 vowels
3 glottis
4 vocal tract

می‌شود که تعداد کلمات محدود باشد و بطور مجزا تلفظ شوند. تشخیص مرز کلمات باعث می‌شود تا در بسیاری از کاربردها مانند بازشناسی گفتار، تبدیل متن به گفتار و ... حجم کمتری از اطلاعات وجود داشته باشد و سرعت افزایش می‌یابد. همچنین قسمت‌های مفید و اصلی گفتار در اختیار می‌باشد که اطلاعات اضافی آن حذف شده در نتیجه دقت نیز افزایش می‌یابد.

۱-۶- روش انجام کار

در این پایان‌نامه از ۹ روش استخراج ویژگی برای تشخیص مرز کلمات استفاده می‌شود. از دو پایگاه داده که دارای سرعت متفاوتی در بیان کلمات هستند، استفاده می‌شود تا تأثیر سرعت تلفظ عبارت گفتاری روی نتیجه نهایی مشخص شود. ابتدا با استفاده از روش مورد نظر ویژگی استخراج می‌شود. سپس برای کاهش پردازش، مقدار DC سیگنال حذف می‌شود. با استفاده از فیلتر باترورث هموارسازی انجام می‌شود و سپس مقدار آستانه اعمال می‌شود. محدوده‌ای از سیگنال که از مقدار آستانه بیشتر باشد به عنوان مرز کلمه در نظر گرفته می‌شود.

۱-۷- ساختار پایان‌نامه

این پایان‌نامه از پنج فصل تشکیل شده است که به شرح زیر می‌باشد:

فصل اول مقدمه‌ای در مورد پردازش گفتار، تشخیص مرز کلمه، استخراج ویژگی و قطعه بندی سیگنال گفتار به بخش‌های آوایی و غیرآوایی می‌باشد و توضیحاتی راجع به هر کدام ارائه شده است. در فصل دوم روش‌های متداولی که برای تشخیص مرز کلمه در کارهای پیشین استفاده شده است، مورد بررسی قرار گرفته و نتایج آن نمایش داده شده است. در فصل سوم روش‌های استخراج ویژگی با استفاده از آنالیز طیفی و آنالیز زمان-فرکانس در سیگنال-های گفتاری مورد بررسی قرار می‌گیرند.

در فصل چهارم ابتدا پایگاه داده‌ای که در این پایان‌نامه مورد استفاده قرار می‌گیرد، معرفی می‌شود و سپس روش‌های پیشنهادی برای تشخیص مرز کلمات معرفی شده است و همچنین نتایج شبیه سازی

حاصل از استخراج ویژگی‌های مختلف نمایش داده شده و تجزیه و تحلیل و نتیجه‌گیری می‌شود.
در فصل پنجم نتیجه‌گیری کلی آورده شده و همچنین برای ادامه کار روش‌هایی پیشنهاد می‌شود.

فصل دوم

مروری بر تحقیقات

پیشین

۲-۱- مقدمه

در فصل قبل راجع به تشخیص مرز کلمه و کاربردهای آن صحبت شد. تشخیص دقیق مرزهای آغاز و پایان کلمه، کیفیت پردازش اولیه را بهبود می بخشد و به کاهش هزینه‌های محاسباتی برنامه‌های گفتاری کمک می‌کند. امروزه مسئله تقسیم‌بندی گفتار و تشخیص نقطه آغازین و پایانی کلمه به کمک روش‌های زمانی و فرکانسی انجام می‌شود. در این فصل روش‌هایی که در کارهای انجام شده برای تشخیص مرز کلمه استفاده شده‌اند، تحلیل و بررسی شده و توضیحاتی راجع به هرکدام از آن‌ها داده می‌شود.

۲-۲- انرژی زمان-کوتاه^۱

انرژی یکی از ویژگی‌های صوتی در گفتار و یک پارامتر مهم در شناسایی مرزها و فعالیت واکداری می‌باشد. در میان نویز، بخش واکدار و بخش بیواک، نویز کمترین انرژی را دارد. این روش به طور عمده در تشخیص مناطق آوایی مورد استفاده قرار می‌گیرد. این روش توانایی تشخیص وجود صدا در محیط‌هایی با SNR پایین را دارد [۱۰]. انرژی زمان-کوتاه، انرژی بخش کوتاهی از سیگنال گفتار است. انرژی زمان-کوتاه یک پارامتر برای طبقه بندی ساده و مؤثر بخش‌های صوتی و بدون صدا است. تعریف انرژی طولانی مدت سیگنال به شرح زیر است [۱۱]:

$$E = \sum_{m=-\infty}^{\infty} x^2(m) \quad (1-2)$$

در عبارت فوق E انرژی سیگنال $x(m)$ را نشان می‌دهد. این تعریف برای سیگنال‌های متغیر با زمان مانند گفتار مناسب نیست. چون گفتار یک سیگنال غیرایستا^۲ است و مشخصات آن با زمان تغییر می‌کند. پس سیگنال گفتار پنجره‌گذاری شده، به فریم‌های کوتاه مدت تقسیم‌بندی می‌شود تا مشخصات سیگنال به سیگنال ایستا^۳ نزدیک شود.

¹ Short-Time Energy (STE)

² non-stationary

³ Stationary

$$E_n = \sum_{m=n-N+1}^n x^2(m) = x^2(n-N+1) + \dots + x^2(n) \quad (2-2)$$

$$x_n(m) = x(m)w(n-m) \quad n-N+1 \leq m \leq n \quad (3-2)$$

انرژی کوتاه مدت از رابطه (۲-۴) به دست می‌آید:

$$E_n = \sum_{m=n-N+1}^n [x(m)w(n-m)]^2 \quad (4-2)$$

$w(n-m)$ پنجره، و n شماره فریم است.

۲-۲-۱- الگوریتم کلی از انرژی زمان-کوتاه

تغییراتی در انرژی مرزهای کلمات مشاهده می‌شود که به عنوان نشانه‌ای برای شناسایی مرزهای کلمه است. در مرزهای کلمه، دامنه به طور ناگهانی بالا یا پایین می‌رود [۱۲]. برای هر سیگنال گفتاری یک آستانه مشخص می‌شود که برای تعیین مرزهای کلمه استفاده می‌شود [۶]. برای مشخص کردن مرزهای کلمه، ابتدا سیگنال گفتار فریم بندی و پنجره‌گذاری می‌شود و سپس انرژی در هر فریم گفتار محاسبه می‌شود. برای محاسبه انرژی از رابطه (۲-۴) استفاده می‌شود. در یک پنجره هر نمونه به توان دو می‌رسد، سپس با هم جمع می‌شوند. محدوده‌ای از سیگنال که انرژی آن از انرژی آستانه بیشتر باشد، به عنوان مرز کلمه در نظر گرفته می‌شود.

۲-۲-۲- بررسی کارهای گذشته

در مرجع [۱۱] انرژی یکی از روش‌هایی است، که برای تشخیص قسمت‌های واکنار و بیواک استفاده شده است. در این آزمایش از دو گوینده مرد و زن استفاده شده است که هر کدام جمله‌ای را به زبان انگلیسی بیان می‌کنند. این پژوهش بیشتر به تأثیر اندازه پنجره پرداخته است. پنجره‌ای که در این آزمایش استفاده شده همینگ می‌باشد. مشاهده شده است هرچه طول پنجره کوچکتر باشد، سیگنال نوسان بیشتری دارد و زمان پردازش نیز کمتر است. ولی اگر طول پنجره بزرگتر باشد سیگنال هموارتر

است و زمان پردازش نیز افزایش پیدا می‌کند. پس انتخاب طول پنجره مناسب، مهم است. شبیه‌سازی با استفاده از نرم‌افزار متلب و ضبط صدا با استفاده از نرم‌افزار پرات^۱ انجام شده است. قسمت‌هایی که انرژی بیشتر است نشان‌دهنده قسمت‌های واگذار بیشتر می‌باشد و قسمت‌های بیواک انرژی کمتری دارند.

در مرجع [۱۰] نیز از انرژی برای تشخیص نقطه پایانی استفاده شده است. سیگنال‌های گفتاری ضبط شده افراد مختلفی هستند که به زبان ماندارین^۲ صحبت می‌کنند. این سیگنال دارای ۲۸ کلمه است و هر کلمه هفت بار گفته می‌شود، پس تعداد کل کلمات ۱۹۶ کلمه می‌باشد. تعداد تشخیص اشتباه با استفاده از روش انرژی و نرخ عبور از صفر ۱۸ مورد است. پس به میزان ۹۰/۸٪ درست تشخیص داده شده است.

۲-۳- نرخ عبور از صفر^۳

یکی از الگوریتم‌های دیگری که برای تشخیص مرزها استفاده می‌شود نرخ عبور از صفر می‌باشد. نرخ عبور از صفر یعنی تعداد دفعاتی که سیگنال گفتار از صفر عبور می‌کند و تغییر علامت می‌دهد. این معیار یک اندازه‌گیری ساده از محتوای فرکانسی یک سیگنال می‌باشد. نرخ عبور از صفر سیگنال گفتار در رابطه (۵-۲) نشان داده می‌شود [۱۱]:

(۵-۲)

$$(۶-۲) Z_n = \sum_{-\infty}^{\infty} |\operatorname{sgn}[x(m)] - \operatorname{sgn}[x(m-1)]| w(n-m)$$

$$\operatorname{sgn}[x(n)] = \begin{cases} 1 & x(n) \geq 0 \\ -1 & x(n) < 0 \end{cases}$$

(۷-۲)

^۱ Praat

^۲ mandarin

^۳ Zero Crossing Rate

$$w(n) = \begin{cases} 1/2(2N) & 0 \leq n \leq N-1 \\ 0 & \text{otherwise} \end{cases}$$

استفاده و N طول پنجره می‌باشد.

۲-۳-۱- نشانه تشخیص مرزها با استفاده از نرخ عبور از صفر

نرخ عبور از صفر برای تشخیص گفتار از نویز و تعیین قسمت شروع و پایان گفتار بسیار مفید است. در حالت سکوت نرخ عبور از صفر سیگنال صفر است. در حالت نویزی نرخ عبور از صفر نسبت به بخش‌های گفتاری بیشتر می‌باشد. با توجه به رابطه (۲-۵) تعداد تغییر وضعیت از مثبت به منفی یا برعکس را شمرده و آن مقدار را به عنوان ZCR در هر فریم در نظر گرفته می‌شود. محدوده‌ای از نمونه‌ها که نرخ عبور از صفر آن‌ها از مقدار آستانه بیشتر است، مرز کلمه در نظر گرفته می‌شود.

۲-۳-۲- بررسی کارهای گذشته

در مقاله [۱۱] علاوه بر انرژی، از روش نرخ عبور از صفر نیز برای تشخیص قسمت‌های واکنار و بیواک استفاده شده است. همانطور که گفته شد در این پژوهش بیشتر به تأثیر اندازه پنجره پرداخته شده است. پنجره مورد استفاده برای بررسی این معیار نیز همین‌گ می‌باشد. مانند روش انرژی مشاهده شده است هرچه طول پنجره کوچکتر باشد، سیگنال نوسان بیشتری دارد ولی اگر طول پنجره بزرگتر باشد سیگنال هموارتر است. طول پنجره تأثیر مهمی در نتیجه نهایی دارد.

در مرجع [۱۳] نرخ عبور از صفر یکی از روش‌هایی است که با روش‌هایی دیگر از جمله انرژی آنتروپی و خطای پیش‌بینی خطی^۱ ترکیب شده و باعث بهبود پاسخ سیستم برای تشخیص مرزها می‌شود. سیستم پیشنهادی از ۱۰ گوینده (۵ مرد و ۵ زن) و ۷۰ عبارت گفتاری از ۷ متن مجزا استفاده کرده است. جدول (۲-۱)، نتایج حداقل دقت به دست و جدول (۲-۲)، نتایج حداکثر دقت بدست آمده از این ارزیابی‌ها را نشان می‌دهد.

¹ Linear Prediction Error

جدول ۱-۲: مقایسه حداقل دقت به دست آمده در سه الگوریتم [۱۳]

Features	Minimum Accuracy
ZCR + Energy + LPE + Spectral Entropy (Proposed Algorithm)	90.35%
ZCR + Energy + Spectral Entropy	89.33%
ZCR + Energy + LPE	89.58%

جدول ۲-۲: مقایسه حداکثر دقت بدست آمده در سه الگوریتم [۱۳]

Features	Maximum Accuracy
ZCR + Energy + LPE + Spectral Entropy(Proposed Algorithm)	97.22%
ZCR + Energy + Spectral Entropy	95.44%
ZCR + Energy + LPE	96.20%

همانطور که در جدول (۱-۲) و (۲-۲) دیده می‌شود، روش پیشنهاد شده در مرجع [۱۳] با حداکثر دقت ۹۷/۲۲٪ و حداقل دقت ۹۰/۳۵٪ بهترین الگوریتم برای شناسایی مرزها بوده است. روش ZRC در ترکیب با سایر روش‌ها باعث بهبود نتایج می‌شود. برخی از خطاهایی که در سیستم رخ می‌دهد، به دلیل کاهش انرژی بخش گفتار و نویز در سیگنال گفتار است.

۲-۴- گام ۱

فرکانس ارتعاش تارهای صوتی فرکانس پایه^۲ یا گام گفته می‌شود و آن را با F_0 نشان می‌دهند. فرکانس پایه تنها برای صداهای آوایی وجود دارد و گام از روش‌هایی است که برای تشخیص مرز کلمه استفاده می‌شود. هدف یافتن دوره تناوب ارتعاش تارهای صوتی در هر فریم از فایل گفتار است. محدوده فرکانس گام، در مردها در حدود ۵۰ تا ۲۵۰ هرتز، در زنان این فرکانس حدود ۱۵۰ تا ۴۵۰ هرتز و در کودکان بین ۳۰۰ تا ۷۰۰ هرتز است. تشخیص گام نقش بسیار مهمی در فشرده‌سازی گفتار، کدگذاری گفتار، شناسایی گوینده، تشخیص و ترکیب گفتار دارد. الگوریتم‌های تشخیص گام زیادی براساس طیف، کپستروم، شبکه‌های عصبی و تبدیل موجک وجود دارد، اما تعداد کمی از روش‌ها مانند تابع میانگین تفاوت اندازه (AMDF)^۳ و تابع همبستگی خودکار (ACF)^۴ در تشخیص پیچ کاربرد دارد. AMDF با گرفتن قدر مطلق بین اختلاف سیگنال گفتار اصلی و سیگنال گفتار شیفت یافته، در هر مقدار تأخیر محاسبه می‌شود. پیاده‌سازی AMDF محاسبات نسبتاً پایینی دارد و پیاده‌سازی آن آسان است. رابطه (۲-۸) تابع میانگین تفاوت اندازه را نشان می‌دهد.

$$x_w(m) = \frac{1}{N-m-1} \sum_{n=0}^{N-m-1} |s_w(n+m) - s_w(n)| \quad (۲-۸)$$

$s_w(n)$ سیگنال گفتار و N طول فریم سیگنال گفتار است. محدوده m بین ۰ تا N می‌باشد. یکی از معایب عمده AMDF برای تشخیص گام، شدت تغییرات سیگنال گفتار و نویز پس زمینه است. به منظور کاهش خطا، پیشنهاد شده است که همراه با AMDF از ACF استفاده شود که ACF در رابطه (۲-۹) نشان داده شده است:

$$R(k) = \sum_{n=0}^{N-k-1} x(n)x(n+k) \quad (۲-۹)$$

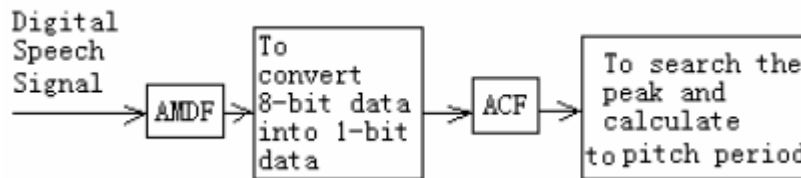
^۱ Pitch

^۲ Fundamental Frequency

^۳ Average Magnitude Difference Function (AMDF)

^۴ Auto Correlation Function (ACF)

ابتدا مقادیر AMDF برای هر فریم سیگنال محاسبه می شود. مقادیر AMDF به سیگنال های یک بیتی تبدیل می شوند تا از پیچیدگی محاسبات در ACF کم شود و همچنین از اثرات دامنه و ساختار سیگنال در تشخیص گام کاسته شوند. سپس مقادیر ACF مقادیر AMDF محاسبه می شود.



شکل ۱-۲: بلوک دیاگرام تشخیص گام با استفاده از AMDF و ACF [۱۴]

میزان خطا با استفاده از الگوریتم تشخیص گام بر اساس AMDF و ACF در مقایسه با AMDF کاهش می یابد، که در جدول (۳-۲) مقایسه شده اند [۱۴].

جدول ۳-۲: مقایسه خطای سه روش برای تشخیص گام [۱۴]

AMDF	ACF	ترکیب هر دو روش
۱۵٪	۹٪	۸٪

۲-۴-۱- الگوریتم تشخیص مرز کلمه با استفاده از روش گام

ابتدا با استفاده از روش های موجود برای تشخیص گام، گام هر فریم محاسبه می شود (در این پایان نامه از روش خودهمبستگی^۱ استفاده شده است). سپس محدوده ای از سیگنال که از مقدار آستانه بیشتر است به عنوان مرز کلمات انتخاب می شود.

۲-۴-۲- انجام کارهای گذشته

یکی از مراجع هایی که از روش گام برای شناسایی مرزهای کلمه استفاده کرده است، مرجع [۱۵] می باشد که از زبان های هندی و آلمانی به عنوان پایگاه داده استفاده کرده است. در زبان هندی مطالعات بر روی گفته های ده گوینده انجام شده و هر گوینده ۵۰ جمله را بیان کرده است. جملات دارای طول

¹ autocorrelation

متفاوت از ۵ کلمه تا ۱۵ کلمه بودند. برای زبان آلمانی، در مجموع از پنج گوینده استفاده شده که هر کدام ده جمله بیان کرده‌اند و به طور متوسط ۱۰ کلمه در هر جمله وجود دارد. نتایج مطالعه به وضوح نشان می‌دهد که زبان‌های هندی برخلاف زبان‌های اروپایی دارای عروضی مشابه هستند. در زبان‌های هندی، اشتباه در تشخیص مرزهای کلمات به دلیل وجود جملات مرکب ایجاد شده است. از نتایج می‌توان دریافت که بیش از ۸۵٪ از مرزهای کلمات در زبان‌های هند به درستی شناسایی شده‌اند، و بیش از ۶۵٪ از مرزهای کلمات با استفاده از روش گام در زبان آلمانی شناسایی می‌شوند.

۲-۵- آنترپی^۱

آنترپی فعالیت‌های تصادفی یک سیستم را اندازه‌گیری می‌کند. مشخص شده است که مقدار آنترپی در تمایز بخش‌های گفتار در یک گفتار پیوسته ضبط شده از قسمت‌های غیرگفتاری، به ویژه در محیط‌های پر سر و صدا، بسیار مفید است. وقتی یک سیستم منظم باشد، آنترپی آن پایین‌تر است و برعکس، هرچه آشفته‌تر سیستم باشد، آنترپی بالاتر است. در سال ۱۹۹۸، شن^۲ و همکاران [۱۶] اولین نفراتی بودند که مفهوم آنترپی را در تشخیص نقطه پایانی گفتار استفاده کردند. ابتدا طیف هر فریم با استفاده از تبدیل فوریه به دست می‌آید. از آنجا که این ویژگی فقط با تصادفی بودن انرژی ارتباط دارد و مستقل از دامنه انرژی سیگنال است، می‌تواند از محاسبات زیادی جلوگیری کرده و در برابر نویز قوی باشد.

اگر $x(n)$ سیگنال گفتار و $x_i(m)$ سیگنال گفتار مربوط به فریم i ام بعد از اضافه کردن پنجره باشد، تبدیل فوریه آن از رابطه (۲-۱۰) به دست می‌آید:

$$X_i(k) = \sum_{m=0}^{N-1} x_i(m) \exp(-i 2\pi km / N) \quad (2-10)$$

و انرژی تبدیل فوریه $x_i(m)$ برابر است با:

$$Y_i(k) = |X_i(k)|^2 \quad (2-11)$$

¹ Entropy

² Shen

تابع چگالی احتمال^۱ طیف برای هر فریم از سیگنال گفتار برآورد می‌شود تا بر اساس آن آنتروپی طیفی تعریف و اندازه‌گیری شود. تابع چگالی احتمال برای طیف را می‌توان با نرمال‌سازی مؤلفه‌های فرکانس تخمین زد:

$$p(k, i) = \frac{Y_i(k)}{\sum_{l=0}^{N/2} Y_i(k)} \quad k = 0, 1, \dots, N/2 \quad (12-2)$$

آنتروپی طیفی هر فریم از رابطه (۱۳-۲) به دست می‌آید:

$$H(i) = -\sum_{k=0}^{N/2} p(k, i) \log p(k, i) \quad (13-2)$$

۲-۵-۱- نشانه‌های تشخیص مرز کلمه با استفاده از روش آنتروپی

آنتروپی گفتار به دلیل ویژگی ذاتی طیف‌های گفتاری با نویز متفاوت است. مزیت اصلی آنتروپی کاهش حساسیت آن نسبت به تغییرات دامنه سیگنال گفتار است، در حالی که حساسیت بیشتری نسبت به وجود یا عدم وجود سیگنال نشان می‌دهد [۱۷]. آنتروپی در سیگنال‌های حاوی نویز دارای بیشترین مقدار، در فرکانس خالص مقدار صفر و در سیگنال گفتار حاوی نویز مقداری بین این دو خواهد داشت [۲]. در فرآیند تشخیص نقطه پایانی، مجموع مقادیر آنتروپی طیفی فریم‌ها ابتدا توسط یک فیلتر هموار می‌شود. مقدار آستانه‌ای در نظر گرفته می‌شود. در محدوده‌ای از سیگنال که مقدار آنتروپی از مقدار آستانه بیشتر باشد مرز شناخته می‌شود.

۲-۵-۲- بررسی کارهای گذشته

فی وو^۲ برای بهبود آنتروپی طیفی، آنتروپی طیفی زیر باند را پیشنهاد کرد و دریافت که این پارامتر مرز بین گفتار و نویز را به وضوح در SNR های ضعیف تقویت می‌کند و برای محیط‌های پر سر و صدا قوی‌تر است. ایده آنتروپی طیفی زیر باند این است که یک فریم به تعدادی زیر باند تقسیم می‌شود،

¹ probability density function (pdf)

² Fei Wu

سپس آنتروپی طیفی زیر باند را بدست می آورد، بنابراین مشکل دامنه که تحت تأثیر نویزاست را از بین می برد [۱۸]. هر باند فرکانسی فریم را می توان به M زیر باند تقسیم کرد، سپس انرژی دامنه فرکانس هر زیر باند را محاسبه کرد. انرژی زیر باند m ام از فریم λ_m به شرح است:

$$E_b(m, i) = \sum_{k=(m-1)*4}^{(m-1)*4+3} Y_i(k) \quad 1 \leq m \leq N_b \quad (14-2)$$

وتابع چگالی احتمال هر زیر باند برابر است با:

$$p_b(m, i) = \frac{E_b(m, i)}{\sum_{k=1}^{N_b} E(m, i)} \quad 1 \leq m \leq N_b \quad (15-2)$$

آنتروپی طیفی هر زیرباند از رابطه (۱۶-۲) به دست می آید:

$$H_b(i) = \sum_{m=1}^{N_b} p_b(m, i) \log p_b(m, i) \quad (16-2)$$

روش پیشنهادی برای تقسیم زیرباندها باعث کاهش خطا و بهبود پاسخ شده است.

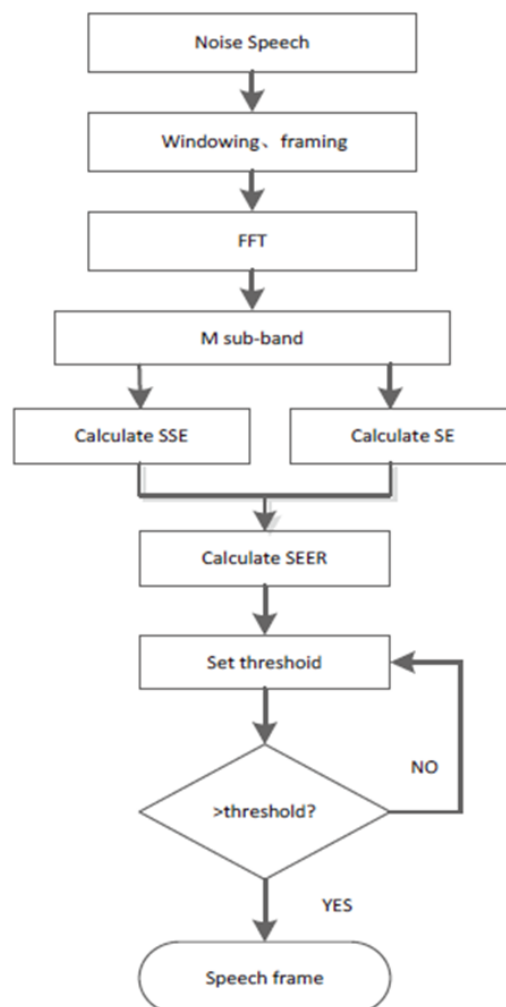
آزمایشات نشان می دهد، انرژی زیر باندها زیاد است و آنتروپی طیفی زیر باند کوچک است. انرژی زیر باند محدوده نویز کوچک است و آنتروپی طیفی زیر باند بزرگ است. پس تفاوت دو محدوده کلمات و نویز باعث می شود تشخیص نقطه پایانی آسان تر باشد. ولی در محیط هایی که نویز زیاد است باعث افزایش اشتباه در تشخیص مرزها می شود. بنابراین در مقاله [۱۹] برای بهبود آنتروپی طیفی زیرباند، الگوریتم تشخیص نقطه پایانی نرخ انرژی-آنتروپی زیرباند^۱ را پیشنهاد می کند که در رابطه (۲-۱۷) نشان داده شده است.

$$SEH(i) = \sqrt{1 + |SE(i) / H_b(i)|} \quad (17-2)$$

ابتدا سیگنال گفتار فریم بندی می شود. در این مقاله از فایل هایی با نرخ نمونه برداری ۸ کیلوهرتز استفاده می کنند، اندازه فریم ۲۵ میلی ثانیه یعنی ۲۰۰ نمونه، شیفت فریم ۱۰ میلی ثانیه، یعنی ۸۰ نمونه است. سپس انرژی زیر باند، آنتروپی طیفی زیر باند و نرخ انرژی-آنتروپی زیرباند محاسبه می شود.

^۱ sub-band energy and sub-band spectral entropy (SEH)

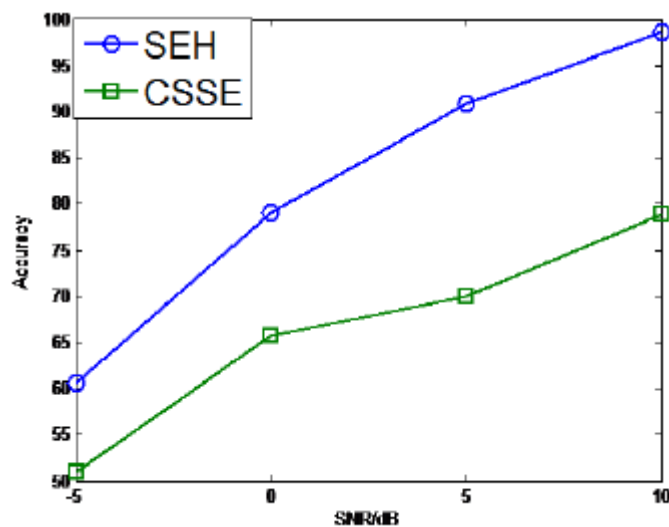
با استفاده از فیلتر برای بهبود سیستم هموار سازی انجام می‌شود. فرایند الگوریتم در شکل (۲-۲) نشان داده شده است. این مقاله با استفاده از Cool Edit Pro پایگاه داده صوتی ایجاد می‌کند، ۵۰ نمونه گفتار بدون نویز را برای آزمایش انتخاب می‌کند، سپس با نویزهای مختلف در SNR های متفاوت آزمایش شده است. این صداها از کتابخانه استاندارد نویز NOISEX92 گرفته شده است. نتایج آزمایش های شبیه سازی به مقایسه دقت نقطه پایانی گفتار آنالوژی طیفی زیر باند معمولی^۱ و الگوریتم بهبود یافته تحت شرایط نویز مختلف و SNR های مختلف می‌پردازد. سپس نمونه‌های گفتار به ترتیب با الگوریتم ارائه شده در این مقاله به تشخیص نقطه پایانی می‌پردازند.



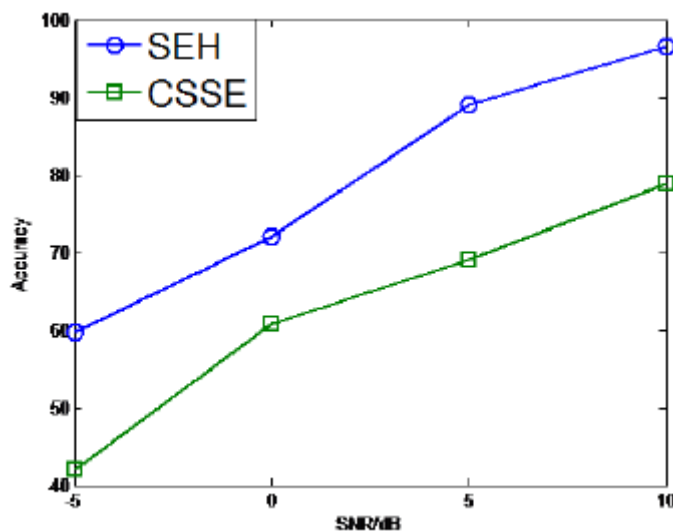
شکل ۲-۲: الگوریتم تشخیص گفتار با استفاده از روش آنالوژی و انرژی زیرباند [۱۹]

¹ conventional sub-band spectral entropy (CSSE)

از شکل (۳-۲) و (۴-۲) مشاهده می‌شود که الگوریتم پیشنهاد شده نسبت به روش آنالیز طیفی تحت SNR های مختلف برتر است. روش آنالیز طیفی مرسوم در حین تشخیص گفتار در SNR های کم دقت کمتری دارد ولی روش بهبود یافته آنالیز طیفی می‌تواند نقاط پایانی را در SNR های بالاتر تشخیص دهد و برای بدست آوردن بخش گفتاری از نویز مفیدتر است، در نتیجه، دقت تشخیص افزایش می‌یابد.



شکل ۳-۲: تشخیص مرزها تحت نویز سفید [۱۹]



شکل ۴-۲: تشخیص مرزها تحت نویز همبسته [۱۹]

فصل سوم

استخراج ویژگی

۳-۱- مقدمه

در فصل قبل روش‌هایی که تاکنون برای تشخیص مرز کلمه استفاده شدند، بررسی شدند. سیگنالی که اطلاعات فرکانسی آن با گذشت زمان تغییر نمی‌کند سیگنال ایستا نامیده می‌شود. در مقابل سیگنال غیرایستا سیگنالی است که اطلاعات فرکانسیش با زمان تغییر می‌کند. نمایش سیگنال در حوزه زمان اطلاعاتی از فرکانس ندارد و زمان اتفاق افتادن فرکانس را نشان نمی‌دهد. همچنین تبدیل فوریه یک سیگنال فقط دارای محتوای فرکانسی است و هیچ اطلاعات زمانی‌ای ندارد. در نمایش سیگنال در حوزه زمان و فرکانس، متغیرهای زمان و فرکانس به صورت مجزا مورد بررسی قرار می‌گیرند. سیگنال گفتار یک سیگنال غیرایستا است و برای تحلیل بهتر باید در حوزه زمان-فرکانس^۱ بررسی شود. نمایش زمان-فرکانس یعنی نمایشی از سیگنال که در هر دو حوزه زمان و فرکانس می‌باشد. نمایش سیگنال در حوزه فرکانس به دلیل عدم توانایی نمایش تغییرات فرکانسی در طول زمان، برای تحلیل سیگنال‌های غیرایستا مناسب نیست. پس نیاز به یک توزیع با توانایی بازنمایی تغییرات فرکانسی نسبت به زمان وجود دارد. توزیع‌های زمان-فرکانس به خوبی این تغییرات سیگنال‌ها را نشان می‌دهند. نمایش زمان-فرکانس هم اطلاعات زمانی دارد و هم اطلاعات فرکانسی و برای سیگنال‌هایی که دارای چندین فرکانس متغیر در زمان هستند، مفید هستند [۲۰]. به طور کلی نمایش زمان-فرکانس یک ابزار قوی جهت پردازش سیگنال‌ها است. در ادامه برخی از روش‌های زمان-فرکانس و روش‌های طیفی بررسی می‌شود.

۳-۲- مقیاس مل^۲

یکی از مقیاس‌های محبوب فرکانس ادراکی، مقیاس مل (مخفف ملودی^۳) است. این مقیاس برگرفته از اثر استونز و همکاران در سال ۱۹۳۷ می‌باشد [۲۱]. شنوایی انسان به اندازه برابر نسبت به همه باندهای فرکانسی یکسان نیست. عملکرد گوش برای فرکانس‌های کمتر از یک کیلوهرتز، خطی و برای

^۱ Time-frequency

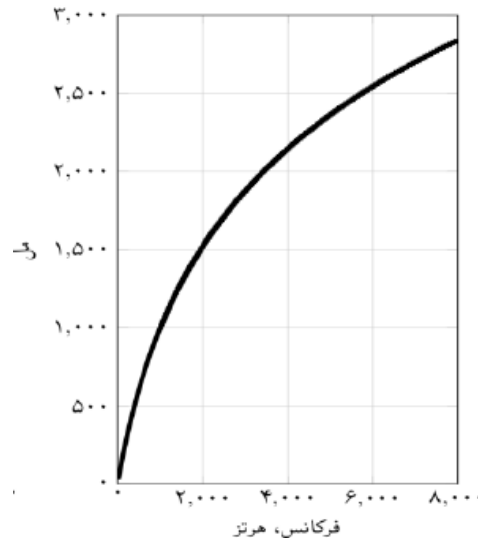
^۲ Mel-scale

^۳ Melody

فرکانس های بالاتر لگاریتمی است. برای تبدیل بین مقیاس ادراکی مل و فرکانس رابطه های متعددی استفاده می شود، اما رابطه ی رایج برای تبدیل فرکانس به مقیاس مل عبارت است از:

$$m(f) = 2595 \times \log_{10} \left(1 + \frac{f}{700} \right) \quad (1-3)$$

شکل (1-3) با توجه به رابطه (1-3)، مقادیر مقیاس مل نسبت به فرکانس را نشان می دهد.



شکل 1-3: تغییرات مقیاس مل برحسب رابطه (1-3) [25]

3-2-1- استخراج ویژگی ضرایب کپسترال در مقیاس مل¹

پردازش MFCC به عنوان پرکاربردترین روش استخراج ویژگی در سیستم تشخیص خودکار گفتار کاربرد دارد که در سال ۱۹۸۰ توسط دیویس^۲ و مرملستین^۳ ابداع شد [21]. MFCC نوعی کپستروم حقیقی است که رفتار سیستم شنوایی انسان را شبیه سازی می کند. مراحل به دست آوردن ضرایب در شکل (2-3) نمایش داده شده است [22].

گوش انسان به اطلاعات فرکانس پایین ارزش بیشتری می دهد. روش MFCC نیز نواحی فرکانس پایین را بهتر از نواحی فرکانس بالا نشان می دهد [23]. بدترین ایراد این روش این است که در مورد

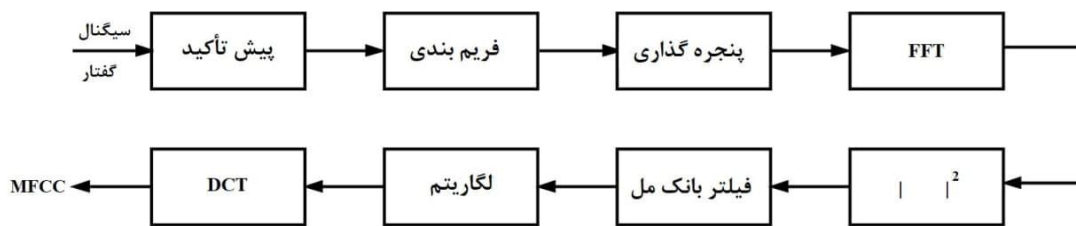
¹ Mel Frequency Cepstral Coefficients (MFCC)

² Davis

³ Mermelstein

نویز ایمن نیست. بهتر از روش انرژی عمل می‌کند ولی در واقع نتایج MFCC در حضور نویز پس زمینه، خیلی قوی نمی‌باشد [۲۴] و [۲۵].

مراحل به دست آوردن ضرایب در شکل (۲-۳) نمایش داده شده است، که هر مرحله توضیح داده می‌شود:



شکل ۲-۳: مراحل استخراج ویژگی MFCC [۲۲]

• پیش تأکید^۱:

در مرحله پیش تأکید از فیلتر بالاگذر برای تقویت فرکانس‌های بالا و حذف اثرات طیفی حنجره و لب‌ها استفاده می‌شود [۲۴].

$$H(Z) = 1 - \mu Z^{-1} \quad (۲-۳)$$

μ پارامتر پیش تأکید و مقدار آن بین ۰/۹ تا ۱ است.

• فریم بندی^۲:

پس از پیش تأکید سیگنال، سیگنال را فریم بندی می‌شود. طول فریم‌ها برای کاربردهای مختلف متفاوت است (بین ۱۰ تا ۵۰ میلی ثانیه). این فریم‌ها می‌توانند همپوشانی^۳ داشته باشند، که این همپوشانی بین ۲۵ درصد تا ۷۵ درصد طول فریم است.

• پنجره گذاری^۴:

برای تبدیل سیگنال غیرایستا به ایستا سیگنال فریم بندی شده در یک پنجره ضرب می‌شود، که

¹ pre-emphasis
² Frame Blocking
³ overlap
⁴ Windowing

معمولا این پنجره، همینگ می باشد. بعد از استفاده از پنجره ناپیوستگی بین مرز فریم‌ها از بین برود.

• آنالیز طیفی:

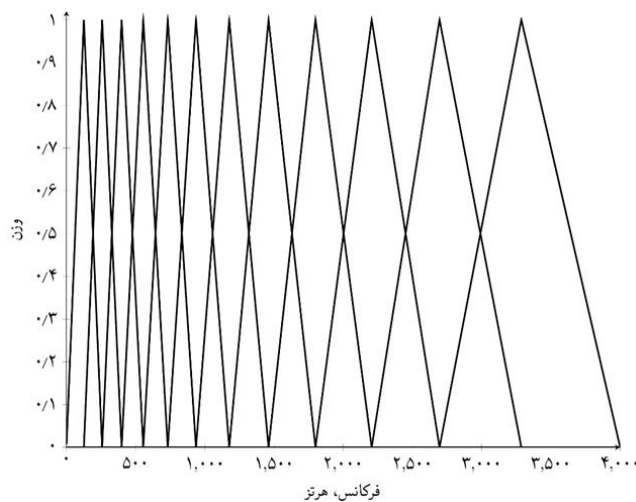
برای انتقال هر فریم گفتار به حوزه فرکانس از تبدیل فوریه گسسته^۱ استفاده می شود آنالیز طیفی نام دارد. چون سیگنال حقیقی و زوج است، به دلیل تقارن نصف سیگنال نگه داشته می شود.

• فیلتر بانک مل:

در مرحله بعد، پردازش فیلتر بانک^۲ مل انجام می شود. فیلتر بانک در مقیاس مل در واقع براساس شنیداری و مدل کردن حساسیت گوش انسان نسبت به حوزه های مختلف فرکانس است که این فرکانس ها با استفاده از رابطه (۱-۳) منتقل می شود. شکل (۳-۳) فیلتر بانک های مثلثی مل را نشان می دهند و در رایج ترین شکل، مجموعه ای از فیلترهای میان گذر (معمولا ۱۲ تا ۳۰ فیلتر) که دارای فرکانس های مختلف است تا محدوده مورد نظر از طیف سیگنال را پوشش دهد [۲۵].

• لگاریتم و تبدیل کسینوسی گسسته:

سپس لگاریتم انرژی فیلتر بانک ها محاسبه می شود و در نهایت با اعمال تبدیل کسینوسی گسسته^۳ به اندازه تعداد فیلتر بانک ها ضرایب کپیسترال مقیاس مل به دست می آید.



شکل ۳-۳: پاسخ فرکانسی فیلتربانک مل [۲۵]

¹ Discrete Fourier Transform (DFT)

² Filter Bank

³ Discrete Cosine Transform (DCT)

۳-۳- تحلیل پیش‌بینی خطی^۱

تحلیل پیش‌بینی خطی یکی از قدرتمندترین و پرکاربردترین روش‌های تحلیل گفتار است. اهمیت این روش هم در توانایی آن در ارائه برآورد دقیق از پارامترهای گفتار و هم در سرعت نسبی محاسبه آن است [۲۶]. مشخصات پیش‌بینی خطی مزایای زیادی در گفتار و تشخیص خودکار گوینده دارد و این ضرایب به راحتی از طریق سیگنال گفتار تعیین می‌شوند. ضرایب پیش‌بینی کننده، اطلاعاتی در مورد فرکانس فرمنت‌ها، پهنای باند آن‌ها و شکل موج چاکنای را نشان می‌دهد. سرانجام، می‌توان از ضرایب پیش‌بینی برای بهبود روش‌های موجود در تشخیص خودکار گوینده و تشخیص مرزها بر اساس اطلاعات استفاده کرد [۲۷].

LP یک روش صحیح و قابل اعتماد برای تخمین پارامترهای سیستم خطی تولید گفتار برای بخش‌های واکنار می‌باشد. در واقع چون واکنه‌ها دارای قطب هستند، فقط از مدل تمام قطب استفاده می‌شود.

$$H(z) = \frac{S(z)}{E(z)} = \frac{1}{1 - \sum_{k=1}^P a_k z^{-k}} \quad (3-3)$$

$S(z)$ سیگنال گفتار، $E(z)$ سیگنال تحریک^۲ و P تعداد ضرایب پیش‌بینی خطی است. به ضرایب a_k ضرایب پیش‌بینی خطی یا ضرایب LP می‌گویند. ایده اولیه پیش‌بینی خطی بر این اساس است که نمونه گفتار فعلی را می‌توان تقریباً به عنوان ترکیب خطی از نمونه‌های گذشته فرض نمود. با استفاده از فیلتر در رابطه (۳-۳) می‌توان نمونه زمانی لحظه n سیگنال خروجی یعنی $s(n)$ را به دست آورد و با استفاده از ترکیب خطی P نمونه قبلی s را پیشگویی کرد که با استفاده از رابطه (۴-۳) نشان داده می‌شود:

$$\tilde{s}[n] = \sum_{k=1}^P a_k s[n-k] \quad (4-3)$$

¹ Linear predictive analysis

² excitation

تفاوت بین نمونه‌های پیش‌بینی شده و نمونه‌های واقعی را خطای پیشگویی می‌گویند که با رابطه (۵-۳) نشان داده می‌شود:

$$e[n] = s[n] - \tilde{s}[n] = s[n] - \sum_{k=1}^p a_k s[n-k] \quad (۵-۳)$$

ضرایب پیش‌بینی خطی را می‌توان با به حداقل رساندن مجموع تفاوت مربع بین نمونه‌های گفتار واقعی و آن‌هایی که به صورت خطی پیش‌بینی شده‌اند، به دست آورد که با رابطه (۶-۳) تعریف می‌شود:

$$E = \sum_{n=0}^{N-1} e^2[n] = \sum_{n=0}^{N-1} \left(s[n] - \sum_{k=1}^p a_k s[n-k] \right)^2 \quad (۶-۳)$$

منجر به مجموعه‌ای از p معادله با p مجهول می‌شود که در رابطه (۷-۳) نشان داده شده است و می‌توان با روش‌های مختلفی آن را حل کرد.

$$\sum_{k=1}^p a_k \phi(i, k) = \phi(i, 0) \quad i = 1, 2, \dots, p \quad (۷-۳)$$

$\phi(i, k)$ تابع همبستگی سیگنال‌های تأخیر یافته توسط i و k می‌باشد. مقدار a_k به دست آمده فرض می‌شود همان مقدار واقعی است که در مدل تولید گفتار به کار می‌رود [۲۸].

۳-۴- مقیاس بارک^۱

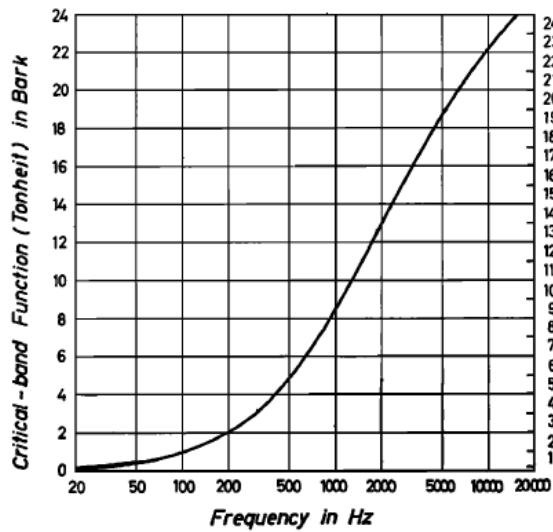
برای تغییر مقیاس فرکانس از مقیاس بارک استفاده شده است که فرکانس صوتی به یک فرکانس غیر خطی تبدیل می‌شود. مقیاس بارک نمایانگر نمایش بهتری از رزلوشن شنوایی انسان است. این مقیاس یک مقیاس آوایی ذهنی است که آقای زویکر^۲ در سال ۱۹۶۰ با رابطه‌ی (۸-۳) نشان داده است [۲۸].

$$Bark(f) = 13a \tan\left(\frac{0.76f}{1000}\right) + 3.5a \tan\left(\frac{f^2}{(7500)^2}\right) \quad (۸-۳)$$

شکل (۵-۳) با توجه به رابطه (۸-۳)، مقادیر مقیاس بارک نسبت به فرکانس را نشان می‌دهد.

^۱ Bark Scale

^۲ Zwicker



شکل ۳-۴: تغییرات مقیاس بارک نسبت به رابطه (۳-۸) [۲۸]

برای مقیاس بارک رابطه‌های متعددی وجود دارد که رابطه‌های (۳-۹) و (۳-۱۰) از رابطه‌های دیگر

این مقیاس می‌باشد.

$$Bark(f) = 25 + 75 \left(1 + 1.4 \left(\frac{f}{1000} \right)^2 \right)^{0.69} \quad (۳-۹)$$

$$Bark(f) = 650 \times \sinh\left(\frac{f}{7}\right) \quad (۳-۱۰)$$

معیار مقیاس مل لگاریتمی و معیار مقیاس بارک عکس تانژانت است.

۳-۴-۱- تحلیل ادراکی پیشگویی خطی^۱

روش PLP تغییر و دگرگونی روش LPC است. از معایب روش LPC این است که مکانیزم سیستم

شنوایی انسان را در نظر نمی‌گیرد و در تمامی فرکانس‌ها سیگنال گفتار را به یک صورت تحلیل می‌کند.

در سال ۱۹۸۹ هرمانسکی^۲ تحلیل ادراکی پیشگویی خطی که مبتنی بر درک انسان است را پیشنهاد

کرد. ایده اصلی روش PLP این است که قدرت تفکیک فرکانسی و حساسیت نسبت به تغییر فرکانس

در گوش انسان در فرکانس‌های مختلف، یکسان نیست و به علاوه حساسیت گوش نسبت به شدت و

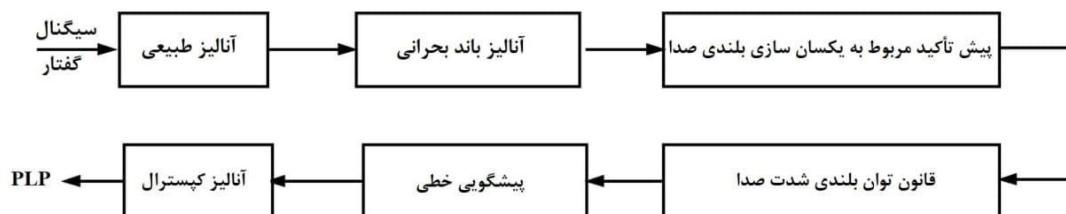
^۱ Perceptual linear predictive (PLP)

^۲ Hermansky

انرژی صوت در فرکانس‌های مختلف متفاوت است [۳۰].

مراحل استخراج ویژگی با استفاده از روش PLP در شکل (۳-۶) آورده شده است که به طور خلاصه

هر محله توضیح داده می‌شود:



شکل ۳-۵: مراحل استخراج ویژگی PLP

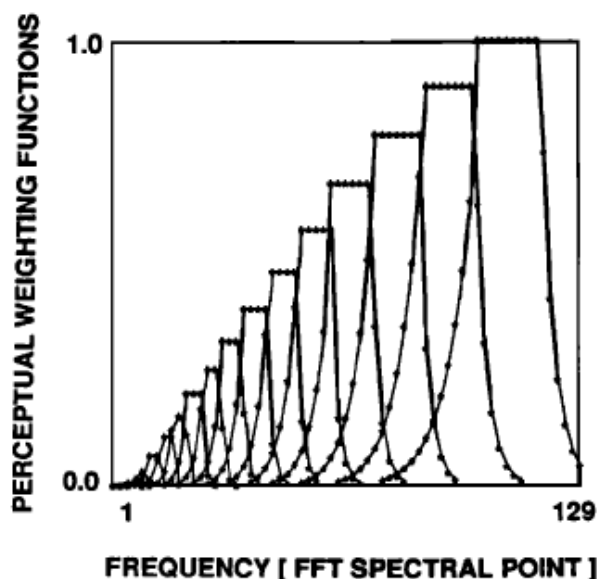
• آنالیز باند طیفی:

مانند روش MFCC ابتدا سیگنال را فریم بندی کرده و سپس در پنجره مناسب ضرب می‌شود. در مرحله آنالیز طیفی طیف توان زمان کوتاه هر سیگنال را با استفاده از تبدیل فوریه گسسته محاسبه می‌شود.

• آنالیز باند بحرانی^۱:

در این مرحله طیف‌های توان به دست آمده با استفاده از تابع انتقال یعنی رابطه (۳-۸) به حوزه بارک برده می‌شوند تا بعضی از فرکانس‌ها بر اساس شنوایی انسان اهمیت بیشتری پیدا کنند در واقع رزولوشن فرکانسی گوش را شبیه‌سازی می‌کند. فیلتربانک‌ها در فرکانس‌های پایین باریک‌تر و در فرکانس‌های بالا وسیع‌ترند و مجموعه‌ای از فیلترهای دوزنقه‌ای به شکل (۳-۶) هستند.

² Critical band analysis



شکل ۳-۶: فیلتربانک بارک [۳۰]

- پیش تأکید مربوط به یکسان سازی بلندی صدا:

این تابع احساس یکسان نبودن بلندی صدا در فرکانس‌های مختلف را جبران می‌کند و رابطه آن به

شکل زیر می‌باشد:

$$E(f) = \left(\frac{f^2}{f^2 + 1.6e^5} \right) \cdot \left(\frac{f^2 + 1.44e^6}{f^2 + 9.6e^6} \right) \quad (۱۱-۳)$$

- قانون توان شدت بلندی صدا^۲:

این تابع رابطه‌ای غیر خطی بین شدت صدا و بلندی صدای دریافت شده را شبیه سازی می‌کند. براین

اساس نشان داده شده است که در گوش انسان میزان احساس بلندی با ریشه سوم انرژی آن متناسب

است در این بلوک رابطه (۱۲-۳) به شکل زیر می‌باشد:

$$\phi(f) = \psi(f)^{0.33} \quad (۱۲-۳)$$

- محاسبه ضرایب پیشگویی خطی:

در این مرحله استفاده از روش LP ضرایب پیشگویی خطی محاسبه می‌شود.

¹ Equal-loudness pre-emphasis

² Intensity to loudness

• آنالیز کیسترال:

بعد از به دست آوردن ضرایب LP، با عکس تبدیل فوریه گرفتن از لگاریتم ضرایب پیشگویی ضرایب PLP محاسبه می‌شود.

برخی آزمایش‌ها نشان می‌دهد که روش MFCC نتایج دقیق‌تری نسبت به روش PLP دارد ولی PLP مقاومت بیشتری نسبت به نویز دارد [۳۱].

۳-۵- تبدیل موجک^۱

تبدیل نگاشتی است که یک فضای ورودی را به یک فضای خروجی می‌برد. هدف این است که اطلاعات بیشتری برای تحلیل سیگنال دریافت شود که به صورت عادی قابل نمایش نیست. تبدیل موجک از مهم‌ترین روش‌های زمان-فرکانسی برای تجزیه و تحلیل سیگنال می‌باشد [۳۲]. موجک با بکارگیری توابعی، بر برخی از محدودیت‌های تحلیل فوریه غلبه می‌کند. در واقع آنالیز فوریه حالت بسیار خاصی از موجک است ولی ضعف فوریه در وهله اول عدم قابلیت نمایش زمان-فرکانس است و نمی‌تواند وضوح^۲ را بالا ببرد و هر جزئیاتی را نمی‌بیند که ناشی از این است که عملکرد آن سرتاسری^۳ است و نمی‌تواند به خوبی وفق^۴ شود در حالی که تبدیل موجک به دلیل خاصیت محلی^۵ بودن می‌تواند به خوبی سیگنال را تحلیل کند. همانطور که گفته شده سیگنال‌های غیرایستا سیگنال‌هایی هستند که فرکانس‌ها و سایر خصوصیات آماری آن‌ها با گذشت زمان تغییر می‌کند. تبدیل موجک برای سیگنال‌های غیرایستا که متناوب نیستند، بسیار مناسب است. تابع موجک قسمتی از شکل را بررسی می‌کند یعنی علاوه بر فرکانس، موقعیت و محل نیز در نظر گرفته می‌شود [۳۳]. یک سیگنال را با استفاده از تبدیل موجک می‌توان به مجموعه‌ای از توابع پایه که به توابع مادر معروف هستند، به صورت پی در پی به

¹ Wavelet Transform

² resolution

³ global

⁴ Adapt

⁵ local

سیگنال‌های تقریب^۱ و جزئیات^۲ تقسیم کرد. اطلاعات فرکانس سیگنال در فرآیند تجزیه به وسیله ضرایب تقریب و ضرایب جزئیات بیان می‌شود. رابطه (۳-۱۳) رابطه بین موجک مادر و فرزندانش می‌باشد:

$$\psi_{a,b}(t) = \frac{1}{\sqrt{a}} \psi\left(\frac{t-b}{a}\right) \quad (13-3)$$

خواص اصلی یک تابع موجک:

(۱) از نظر زمانی، به طرز مناسبی باید محدود باشد. یعنی:

بزرگترین مقدار LB وجود داشته باشد، به نحوی که:

$$\forall t \leq LB : \psi(t) = 0 \quad (14-3)$$

و کوچکترین مقدار UB وجود داشته باشد، به نحوی که:

$$\forall t \geq UB : \psi(t) = 0 \quad (15-3)$$

(۲) باید دارای میانگین صفر باشد:

$$\int_{-\infty}^{\infty} \psi(t) dt = 0 \quad (16-3)$$

(۳) نرم غیر صفر داشته باشد:

$$0 < \int_{-\infty}^{\infty} |\psi(t)|^2 dt < \infty \quad (17-3)$$

صدای انسان از دو مؤلفه فرکانس پایین و فرکانس بالا تشکیل شده است. با اعمال فیلترهای پایین گذر^۳ و بالاگذر^۴ بر روی سیگنال اصلی و در مرحله بعد بر روی سیگنال تقریب، سیگنال‌های تقریب و جزئیات به دست می‌آیند. خروجی فیلتر پایین گذر مؤلفه تقریب و خروجی فیلتر بالاگذر مؤلفه جزئیات است. شکل (۳-۷) تبدیل موجک گسسته^۵ با اعمال فیلتربانک‌ها نشان می‌دهد. پس از اعمال

¹ Approximating

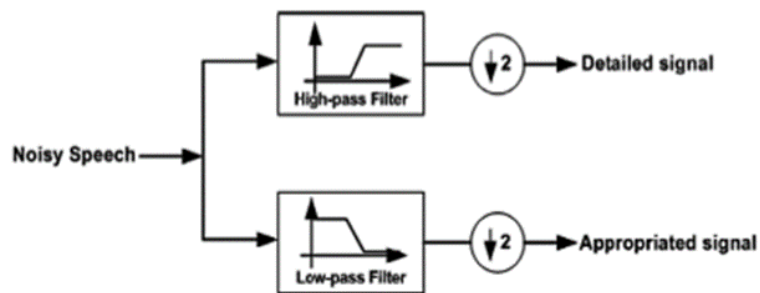
² Detail

³ Low-pass filter

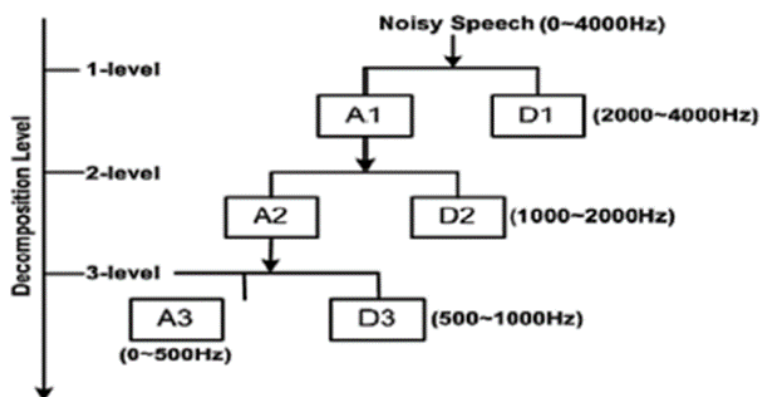
⁴ High-pass filter

⁵ Discrete Wavelet Transform (DWT)

فیلتر، نمونه برداری با فاکتور دو انجام می‌شود و حجم اطلاعات کاهش می‌یابد. آن قدر این کار را ادامه داده می‌شود تا سیگنال تقریباً صفر یا نزدیک صفر شود. تعداد سطح‌ها^۱ می‌تواند متغیر باشد. هرچه سیگنال ساده‌تر باشد مؤلفه تقریباً زودتر به صفر می‌رسد. همچنین می‌توان با اعمال یک بلوک بازسازی کننده^۲ سیگنال اصلی یا تقریبی از آن را بدست آورد. در شکل (۳-۸) نمودار مراحل تجزیه یک سیگنال به سه سطح نشان داده شده است. [۳۴].



شکل ۳-۷: تبدیل موجک گسسته با استفاده از فیلتر بانک [۳۴]

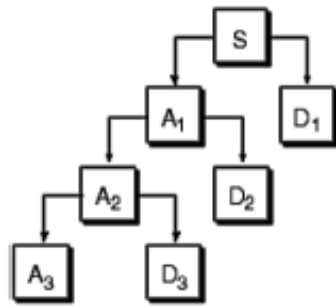


شکل ۳-۸: نمودار مراحل تجزیه سیگنال به سه سطح [۳۴]

تعداد نمونه‌ها پس از تجزیه ثابت می‌باشد که در شکل (۳-۹) نشان داده شده است.

¹ level

² reconstructor



$$\begin{aligned}
 S &= A_1 + D_1 \\
 &= A_2 + D_2 + D_1 \\
 &= A_3 + D_3 + D_2 + D_1
 \end{aligned}$$

شکل ۳-۹: مراحل تجزیه سیگنال به سه سطح [۳۴]

در حقیقت، یک گفتار آوایی دارای انرژی بیشتری نسبت به گفتار غیر آوایی در سیگنال است. بنابراین،

می توان از انرژی سیگنال زیر باند برای استخراج فعالیت گفتاری استفاده کرد [۳۴].

۳-۶- توزیع ویگنر^۱

تابع توزیع ویگنر اولین بار توسط ایگن ویگنر^۲ در سال ۱۹۳۲ برای توضیح مکانیک کوانتوم مورد

استفاده قرار گرفت [۳۵] و در ۱۹۴۸ توسط ویل^۳ در پردازش سیگنال مورد استفاده قرار گرفت [۳۶].

این تابع یکی از توزیع های مهم در حوزه زمان-فرکانس است و جزو توزیع های درجه دوم می باشد. یعنی

توزیعی که یک سیگنال را از حوزه زمان به صفحه چگالی انرژی نگاشت می دهد و این توزیع ها همزمان

تفکیک پذیری فرکانسی و زمانی بالایی دارند [۳۷].

توزیع ویگنر-ویل برای یک سیگنال به صورت تبدیل فوریه تابع همبستگی وابسته به زمان که به

اندازه τ اختلاف زمانی دارند به صورت رابطه (۳-۱۸) می باشد:

$$W_x(t, \nu) = \int_{-\infty}^{+\infty} x(t + \tau/2) x^*(t - \tau/2) e^{-j2\pi\nu\tau} d\tau \quad (3-18)$$

سیگنال چیرپ^۴ در شکل (۳-۱۱) نشان داده شده است. توزیع ویگنر این سیگنال شکل (۳-۱۱) است.

ویگنر یک روش تحلیل سیگنال در حوزه زمان-فرکانس است. مشاهده می شود که شکل (۳-۱۱) هم

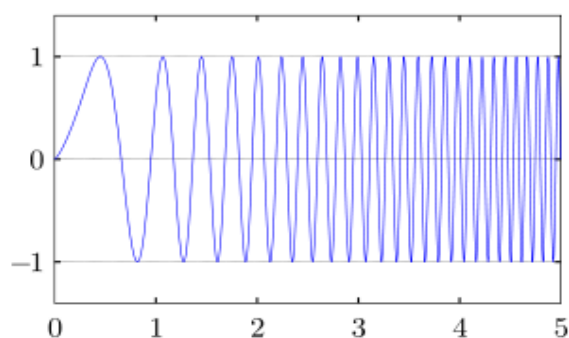
اطلاعات زمانی وهم اطلاعات فرکانسی دارد، با افزایش زمان فرکانس نیز افزایش می یابد.

¹ Wigner Distribution

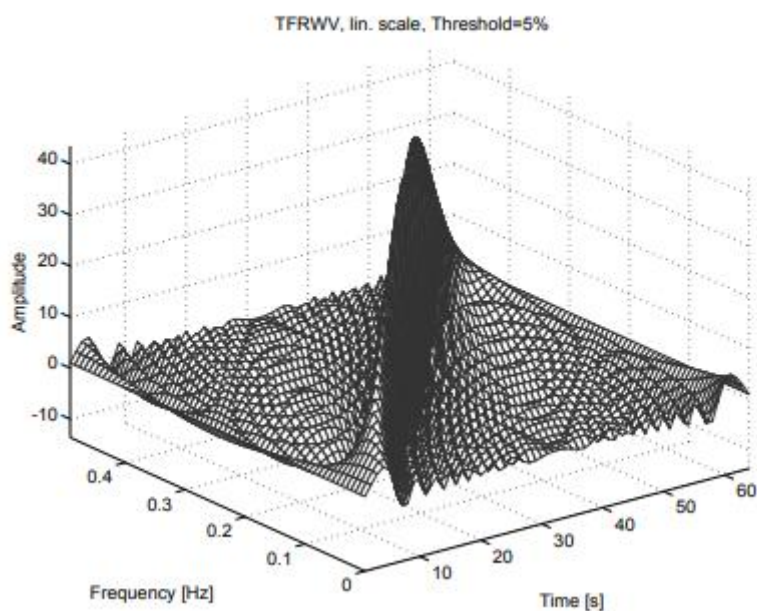
² Eugene Wigner

³ ville

⁴ chirp



شکل ۳-۷: سیگنال چیرپ [۴۰]



شکل ۳-۸: تبدیل ویگنر-ویل سیگنال چیرپ [۴۰]

خصوصیات مختلف و مطلوب توزیع ویگنر مانند تفکیک پذیری زمانی و فرکانسی بالا و موارد دیگر، توزیع ویگنر را به ابزاری مفید برای تجزیه و تحلیل سیگنال تبدیل می‌کند. اشکال اصلی این توزیع درجه دوم بودن آن و ایجاد عبارات متقاطع^۱ بین فرکانس‌های مثبت و منفی سیگنال است که باعث می‌شود تفسیر فضای تبدیل را دشوار می‌کند. از راه‌هایی که می‌توان توزیع ویگنر را بهبود بخشید از

^۱ cross terms

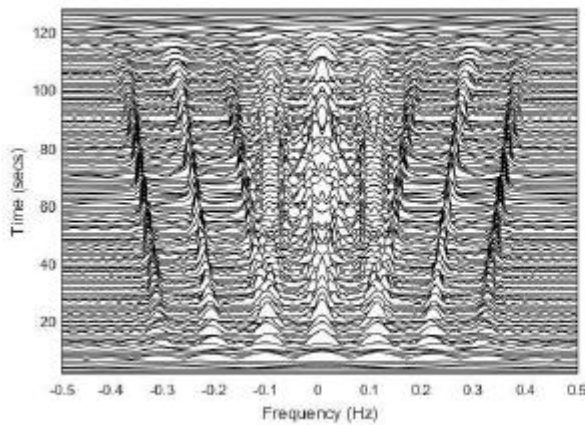
بین بردن فرکانس‌های منفی است تا عبارات متقاطع ایجاد نشود. عبارات متقاطع با مخلوط کردن، فیلتر کردن و یا اضافه کردن تابعی از بین می‌روند. به سیگنالی که در آن عبارات متقاطع وجود ندارد، سیگنال تحلیلی گفته می‌شود. توزیع ویگنر-ویل همان توزیع ویگنر است با این تفاوت که به جای سیگنال اصلی از سیگنال تحلیلی استفاده می‌شود. برای یک سیگنال که از M مؤلفه تشکیل شده است تبدیل ویگنر به صورت زیر می‌باشد [۳۸]:

$$x(t) = \sum_{i=1}^M x_i(t) \quad (۱۹-۳)$$

$$W_x(t, \omega) = \sum_{i=1}^M W_{x_i}(t, \omega) + \sum_{k=1}^{M-1} \sum_{l=k+1}^M 2\Re[W_{x_k x_l}(t, \omega)] \quad (۲۰-۳)$$

در رابطه (۲۰-۳) عبارت دوم جمله تداخلی می‌باشد. در واقع مجموع ویگنر دو سیگنال دقیقاً برابر

با ویگنر مجموع دو سیگنال نمی‌باشد.



شکل ۳-۹: وجود عبارات متقاطع در توزیع ویگنر [۳۷]

در شکل (۳-۱۲) توزیع ویگنر یک سیگنال دو مؤلفه‌ای می‌باشد. بخش‌های ۲ و ۶ به ترتیب عبارات متقاطع بین مؤلفه‌های سیگنال ۱ و ۳ و همچنین ۵ و ۷ هستند. بخش ۴ عبارت متقاطع مضاعف حاصل از مؤلفه‌های ۳ و ۵ و همچنین ۱ و ۷ است. یعنی بخش ۴ دارای دو عبارت متقاطع است. بخش ۸ عبارت متقاطع بین مؤلفه‌های ۱ و ۵ و بخش ۹ عبارت متقاطع بین مؤلفه‌های ۳ و ۷ است. پس عبارت متقاطع بین دو مؤلفه سیگنال اصلی قرار می‌گیرد [۳۷].

۷-۳- ویژگی‌های توزیع ویگنر-ویل

توزیع ویگنر دارای خواص و ویژگی‌هایی می‌باشد که باعث شده این توزیع در برخی از کاربردها مفید باشد. چند نمونه از مهم‌ترین ویژگی‌ها در زیر بیان شده است [۴۰]:

(۱) نگه داری انرژی^۱:

با انتگرال‌گیری از توزیع ویگنر-ویل تابع X روی کل صفحه زمان-فرکانس انرژی کل سیگنال به صورت رابطه نشان (۳-۲۱) داده می‌شود:

$$E_x = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} W_x(t, v) dt dv \quad (۳-۲۱)$$

(۲) شرایط حاشیه‌ای^۲:

چگالی طیفی انرژی و توان لحظه‌ای سیگنال x را می‌توان به عنوان شرایط حاشیه‌ای توزیع ویگنر-ویل از روابط (۳-۲۲) و (۳-۲۳) به دست آورد. شرایط حاشیه‌ای یعنی انتگرال در راستای فرکانس صفحه‌ی زمان-فرکانس که چگالی انرژی است و انتگرال در راستای زمان صفحه‌ی زمان-فرکانس که برابر با طیف توان سیگنال است، وجود داشته باشد.

(۳-۲۲)

$$\int_{-\infty}^{+\infty} W_x(t, v) dt = |X(v)|^2$$
$$\int_{-\infty}^{+\infty} W_x(t, v) dv = |X(t)|^2 \quad (۳-۲۳)$$

(۳) مقادیر حقیقی^۳:

تابع توزیع ویگنر-ویل به ازای هر مقدار t و v دارای مقادیر حقیقی است.

$$W_x(t, v) \in \mathbb{R}, \quad \forall t, v \quad (۳-۲۴)$$

(۴) شیفت دادن:

¹ Energy conservation

² Marginal properties

³ Real-valued

از دیگر خواص توزیع ویگنر-ویل شیفت دادن است. یعنی اگر سیگنال اصلی در زمان یا فرکانس شیفت پیدا کند توزیع ویگنر-ویل نیز در همان زمان یا فرکانس همانند روابط (۲۵-۳) و (۲۶-۳) شیفت پیدا می‌کند.

$$y(t) = x(t - t_0) \Rightarrow W_y(t, \nu) = W_x(t - t_0, \nu) \quad (25-3)$$

$$y(t) = x(t) e^{j2\pi\nu t} \Rightarrow W_y(t, \nu) = W_x(t, \nu - \nu_0) \quad (26-3)$$

(۵) فیلتر کردن:

اگر سیگنال y برابر کانوولوشن x و h باشد، بنابراین توزیع ویگنر-ویل سیگنال y برابر با کانوولوشن زمانی سیگنال‌های h و x است که در رابطه (۲۷-۳) نشان داده شده است.

$$y(t) = \int_{-\infty}^{+\infty} h(t-s) x(s) ds \Rightarrow W_y(t, \nu) = \int_{-\infty}^{+\infty} W_h(t-s, \nu) W_x(s, \nu) ds \quad (27-3)$$

(۶) فرکانس لحظه‌ای^۱:

یعنی فرکانس سیگنال در یک لحظه خاص را نشان می‌دهد و از گشتاور مرتبه اول در فرکانس توزیع ویگنر-ویل به دست می‌آید که در رابطه (۲۸-۳) نشان داده شده است:

$$f_x(t) = \frac{\int_{-\infty}^{+\infty} \nu W_{xa}(t, \nu) d\nu}{\int_{-\infty}^{+\infty} W_{xa}(t, \nu) d\nu} \quad (28-3)$$

(۷) تأخیر گروهی^۲:

این کمیت از گشتاور مرتبه اول در زمان توزیع ویگنر-ویل به دست می‌آید که در رابطه (۲۹-۳) نشان داده شده است.

$$t_x(\nu) = \frac{\int_{-\infty}^{+\infty} t W_{xa}(t, \nu) dt}{\int_{-\infty}^{+\infty} W_{xa}(t, \nu) dt} \quad (29-3)$$

(۸) خاصیت پارسوال:

¹ Instantaneous frequency

² Group delay

اندازه انرژی در حوزه زمان با انرژی در حوزه فرکانس برابر است، که در رابطه (۳-۳۰) نشان داده شده

است:

$$\left| \int_{-\infty}^{+\infty} x(t)y^*(t)dt \right|^2 = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} W_x(t,\nu)W_y^*(t,\nu)dtd\nu \quad (۳-۳۰)$$

فصل چهارم

روش پیشنهادی و نتایج

شبیه سازی

۴-۱- مقدمه

در این فصل ابتدا پایگاه داده مورد استفاده معرفی می‌شود و سپس نتایج شبیه‌سازی الگوریتم‌های استخراج ویژگی گوناگون برای تشخیص مرز کلمه که در این پایان نامه آورده شده است، نمایش داده می‌شود. نتایج تحلیل می‌شوند و بهترین نتیجه مشخص می‌گردد.

۴-۲- پایگاه داده^۱

پایگاه داده TIMIT یکی از معروفترین و پرکاربردترین پایگاه داده انگلیسی مورد استفاده در پردازش گفتار است که نتیجه تلاش مشترک محققان MIT، SRI و TI است. همانطور که در جدول (۴-۱) نشان داده شده است، این پایگاه داده از ۲۳۴۲ جمله مجزا از سه مجموعه مختلف تشکیل شده است. ضبط داده‌ها توسط محققان TI انجام شده است. در مجموع ۶۳۰۰ جمله جمع‌آوری شد که از هر ۶۳۰ گوینده، ده جمله ضبط شده است. تقسیم‌بندی جملات در جدول (۴-۱) نشان داده شده است.

جدول ۴-۱: تقسیم بندی جملات پایگاه داده TIMIT [۴۱]

Sentence Type	Sentences	Talker	Sentences/Talker
Dialect (SIR)	2	1260	2
Compact (MIT)	450	3150	5
Diverse (TI)	1890	1890	3
Total	2342	3600	10

جملات مجموعه Dialect مربوط به لهجه می‌باشند. جملات مجموعه Compact از لحاظ واج متراکم می‌باشند و جملات مجموعه Diverse دارای واج‌های مختلف هستند و تقریباً دارای همه‌ی واج‌ها می‌باشد. این پایگاه داده به دو دسته‌ی آزمایش^۲ و آموزش^۳ تقسیم بندی می‌شود که حدود ۲۰ تا ۳۰ درصد مربوط به آزمایش و حدود ۷۰ تا ۸۰ درصد مربوط به آموزش می‌باشد. هیچ‌کدام از گویندگان در

¹ data base

² test

³ train

هر دو دسته عضویت ندارند. تعداد جملات، تعداد کلمات و بعضی اطلاعات دیگر در جدول (۲-۴) نشان داده شده است [۴۱].

جدول ۲-۴: شرح پایگاه داده TIMIT [۴۱]

تعداد جملات	۵۰۴۰
تعداد کلمات	۴۱۱۶۱
میانگین تعداد کلمات خاص	۵۱۰۷
میانگین تعداد کلمات در هر جمله	۸/۲
میانگین تعداد هجا در هر کلمه	۱/۵

۴-۳- پیش پردازش^۱

به منظور استخراج ویژگی‌ها از سیگنال گفتار، باید ابتدا یک سری پیش‌پردازش‌های اولیه روی سیگنال انجام شود. در این پایان‌نامه برای شبیه‌سازی نتایج الگوریتم‌های مورد نظر از نرم‌افزار متلب استفاده شده است. سعی شده است که موارد مورد استفاده در تمامی الگوریتم‌ها یکسان باشد تا بهتر بتوان نتایج را مقایسه کرد. پنجره مورد استفاده همینگ و طول پنجره ۳۲۰ نمونه است. میزان شیفت پنجره ۲۰ نمونه می‌باشد. همچنین به منظور هموارسازی سیگنال و حذف فرکانس‌های بالا از فیلتر باترورث^۲ استفاده می‌شود. در پردازش سیگنال برای هموارسازی شکل‌ها معمولاً از این فیلتر استفاده می‌شود. برای شبیه‌سازی الگوریتم‌ها از دونوع سیگنال گفتار استفاده شده است:

داده اول مربوط به CADLab است که جملات به طور واضح و شمرده گفته شده است [۴۲]. داده دوم از پایگاه داده TIMIT است که گوینده به طور طبیعی و تقریباً سریع صحبت کرده و و زمان بین کلمات کم است.

¹ Preprocessing

² Butterworth

مقادیر آستانه در هر مرحله به صورت تجربی و با سعی و خطا بدست می‌آید تا نتیجه بهتری حاصل شود. یعنی مقادیرهای مختلفی برای آستانه در نظر گرفته می‌شود و نتیجه شبیه‌سازی با توجه به آن مقادیر آستانه نمایش داده می‌شود. مقدار آستانه‌ای که بهترین نتیجه را ارائه می‌دهد، انتخاب می‌شود. نتیجه شبیه‌سازی روش‌های مختلف برای تشخیص مرز کلمات تصویری است که دارای ۴ بخش می‌باشد. بخش اول مرزهای کلمه روی سیگنال گفتار با استفاده از اطلاعات موجود در پایگاه داده و به صورت دستی مشخص شده است. در بخش دوم با استفاده از روش‌های مختلف ویژگی‌های مورد نظر استخراج می‌شوند. در بخش سوم با استفاده از فیلتر باتروث هموار سازی انجام می‌شود تا تغییرات سیگنال نرم‌تر شود و راحت‌تر بتوان مرزها را تشخیص داد و مقدار آستانه اعمال می‌شود. در بخش آخر مرزهای کلمه با استفاده از روش استخراج ویژگی مورد نظر به دست آمده‌اند. لازم به یادآوری است که دلیل این‌که پایگاه داده مورد استفاده در این تحقیق با پایگاه داده مقالات یکی نبودند، ابتدا روش‌های متداول شبیه‌سازی شده‌اند و نتایج آورده شده است.

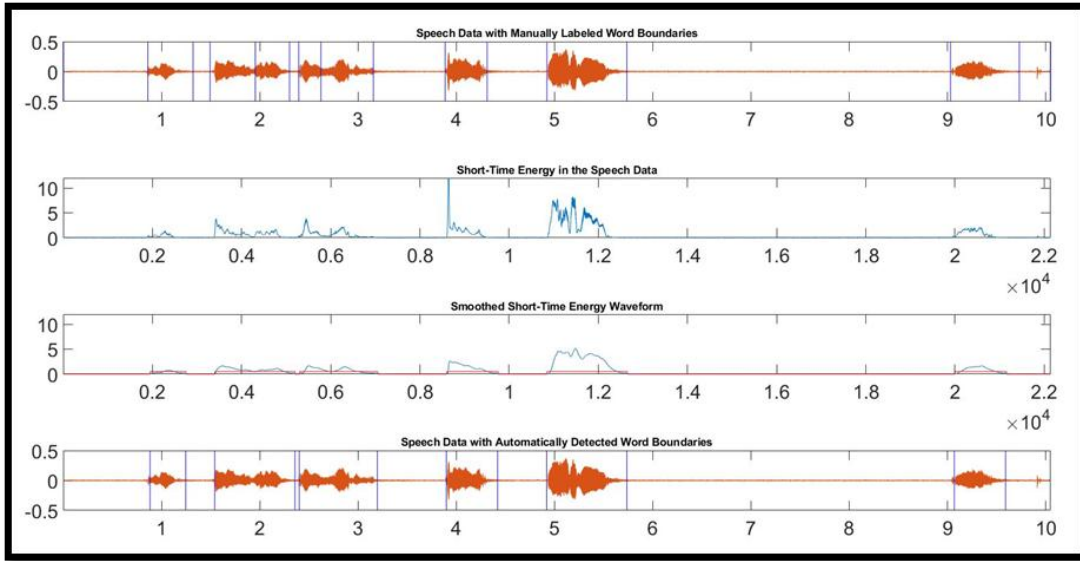
۴-۴- نتایج شبیه‌سازی روش‌های متداول

در ابتدا نتایج روش‌های متداولی که از نوع تحلیل زمانی هستند و برای تشخیص مرزهای کلمه استفاده می‌شدند، بررسی می‌شود.

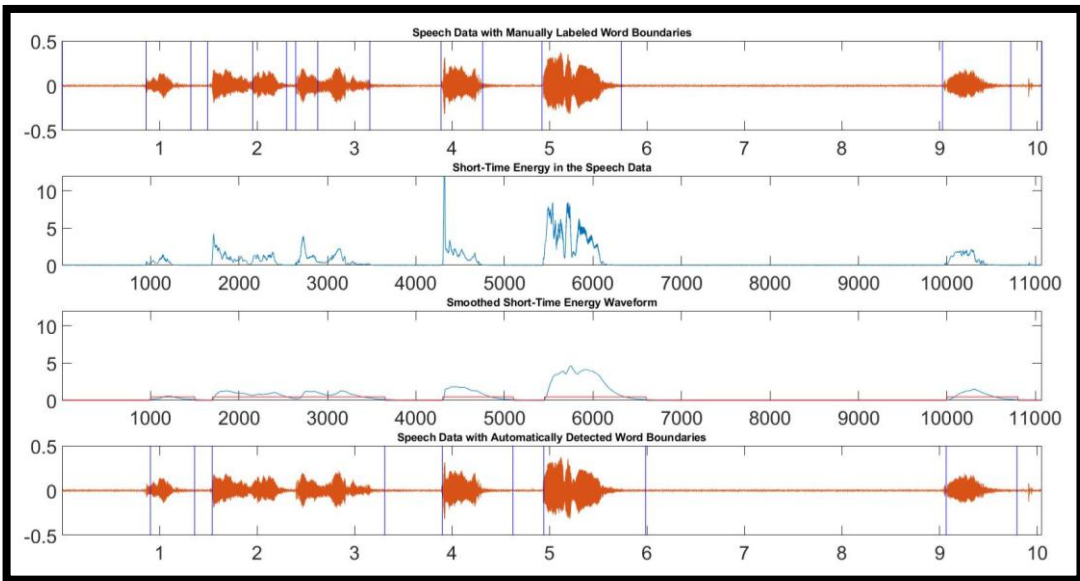
۴-۴-۱- نتایج شبیه‌سازی روش انرژی زمان-کوتاه

الگوریتم جداسازی مرز کلمات با استفاده از انرژی زمان-کوتاه بدین صورت است که ابتدا به منظور بررسی زمان-کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره‌هایی به طول ۳۲۰ نمونه و با شیفتهای زمانی ۲۰، ۴۰ و ۸۰ نمونه، فریم بندی و پنجره‌گذاری می‌گردد و سپس انرژی هر فریم محاسبه می‌شود. برای محاسبه انرژی، نمونه‌های هر پنجره را به توان دو رسانده و باهم جمع می‌شوند. در مرحله بعد به منظور نرم‌تر کردن تغییرات منحنی مربوط به انرژی، مقادیر بدست آمده برای انرژی زمان کوتاه توسط فیلتر باتروث فیلتر می‌گردد. در نهایت با توجه به میزان آستانه تعریف شده، محدوده

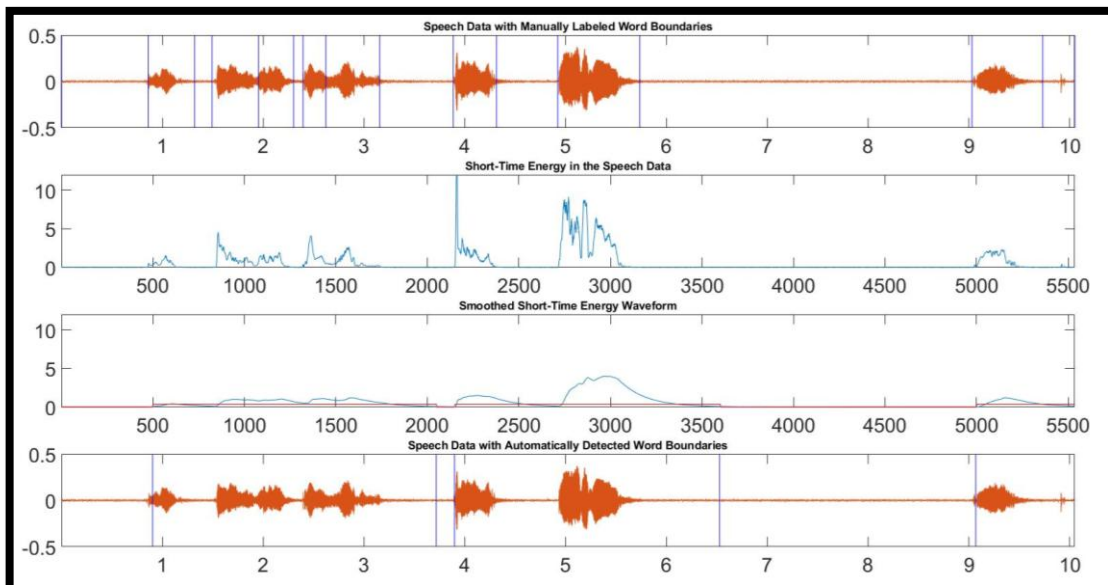
هایی از انرژی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می گردد.



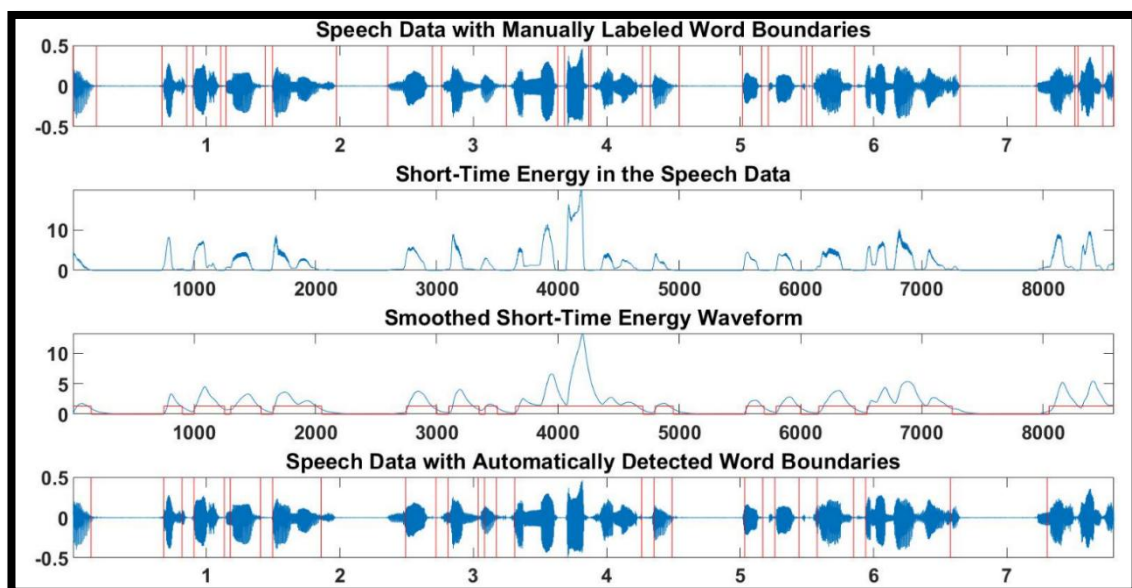
شکل ۱-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه



شکل ۲-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۴۰ نمونه



شکل ۳-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده اول با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۸۰ نمونه



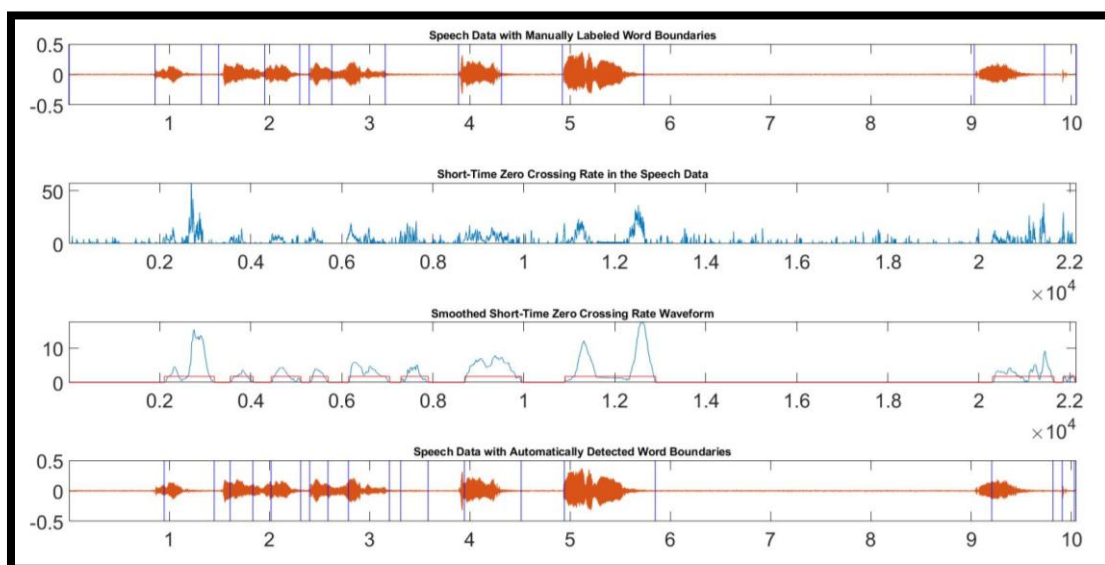
شکل ۴-۴: تشخیص مرز کلمه با استفاده از روش انرژی برای داده دوم با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه

با مقایسه نتایج بدست آمده در اشکال (۱-۴) تا (۴-۴)، مشاهده می‌شود که میزان شیفت زمانی در دقت جداسازی کلمات بسیار تاثیرگذار است. همانطور که در فصل دوم بررسی شد، اگر ابعاد پنجره کوچک باشد سیگنال نوسان زیادی دارد ولی اگر ابعاد پنجره افزایش یابد سیگنال هموارتر است. اگر ابعاد پنجره بدون آن که میزان شیفت زمانی تغییر کند، افزایش یابد، زمان پردازش بیشتر می‌گردد بدون آنکه دقت جداسازی تغییر کند. اگر شیفت زمانی خیلی کوچک انتخاب شود، زمان پردازش به

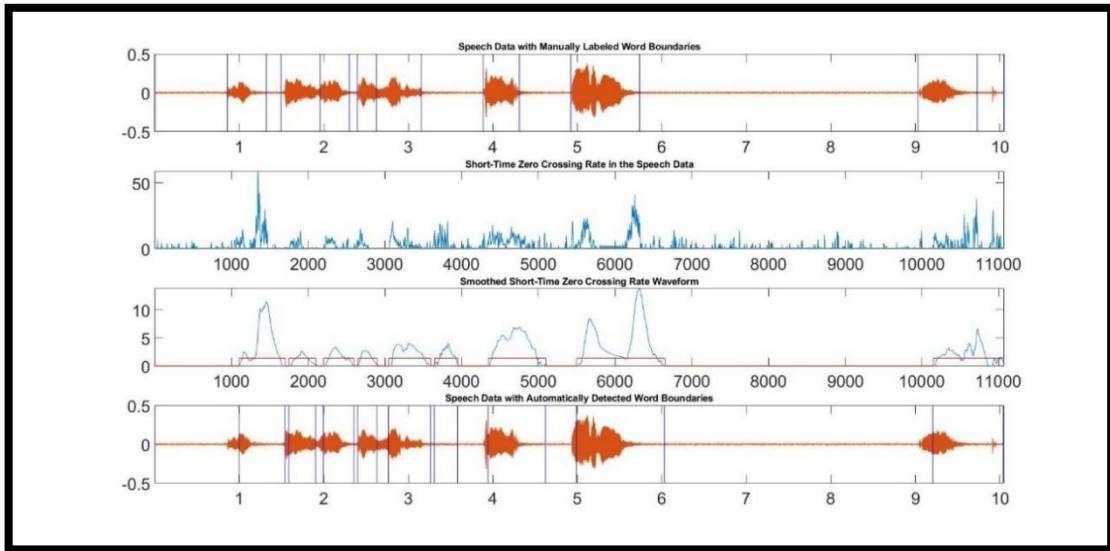
شدت افزایش یافته و دقت جداسازی کلمات نیز کاهش می یابد. بدین صورت که مرزهای بیشتری برای کلمات پیدا می کند. از طرف دیگر اگر میزان شیفت زمانی خیلی بزرگ (به نسبت ابعاد پنجره) انتخاب گردد، در اینصورت دقت جداسازی کلمات کاهش می یابد در این حالت هر چند کلمه را به عنوان یک بخش در نظر می گیرد. به همین خاطر انتخاب ابعاد پنجره و مهم تر از آن میزان شیفت زمانی در دقت جداسازی مرز کلمات تاثیر بسیار زیادی دارد.

۴-۲-۴- نتایج شبیه سازی روش نرخ عبور از صفر

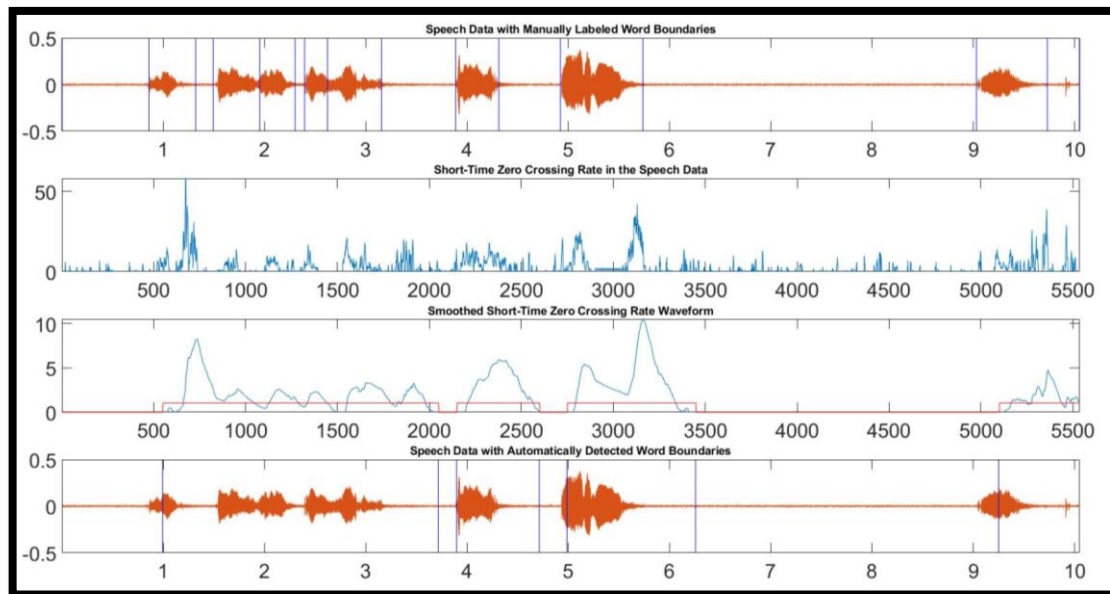
برای تشخیص مرز کلمات با استفاده از معیار گذر از صفر زمان-کوتاه ابتدا به منظور بررسی زمان-کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره هایی به طول ۳۲۰ نمونه و با شیفت های زمانی ۲۰، ۴۰ و ۸۰ نمونه، فریم بندی می گردد و سپس نرخ عبور از صفر هر فریم محاسبه می شود. چون تغییرات خروجی سیگنال خروجی زیاد است، مقدار DC خروجی این مرحله حذف می شود تا حول صفر نوسان کند. در مرحله بعد به منظور نرم تر کردن تغییرات منحنی مربوط به این معیار، مقادیر بدست آمده برای مقادیر گذر از صفر زمان-کوتاه توسط فیلتر باترورث فیلتر می گردد و در نهایت با توجه به میزان آستانه تعریف شده، محدوده هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می گردد.



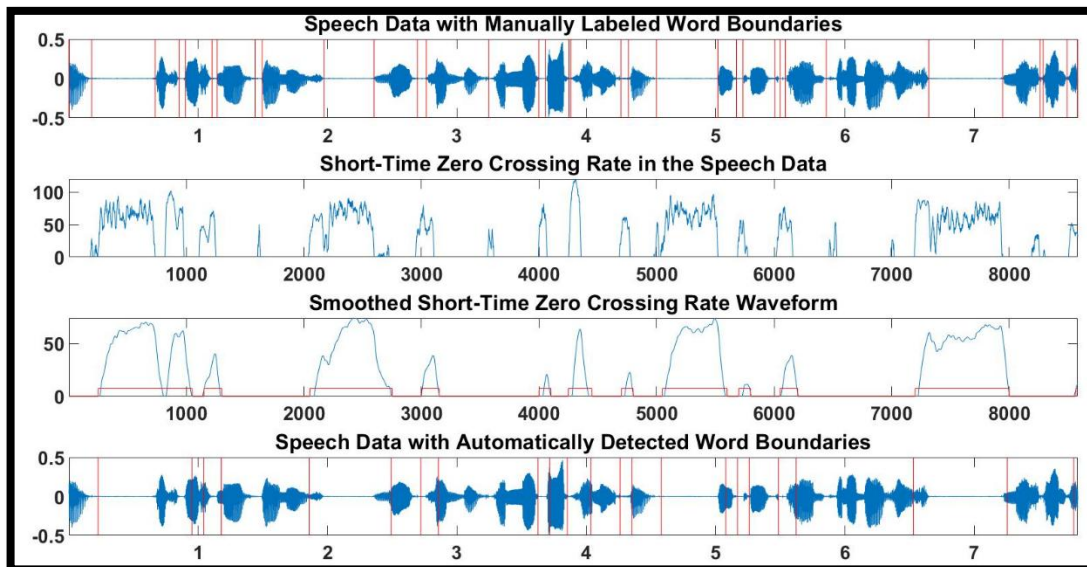
شکل ۴-۵: تشخیص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه



شکل ۴-۶: تشخیص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شیفت زمانی ۴۰ نمونه



شکل ۴-۷: تشخیص مرز کلمه با استفاده از نرخ عبور از صفر برای داده اول و طول پنجره ۳۲۰ نمونه و شیفت زمانی ۸۰ نمونه

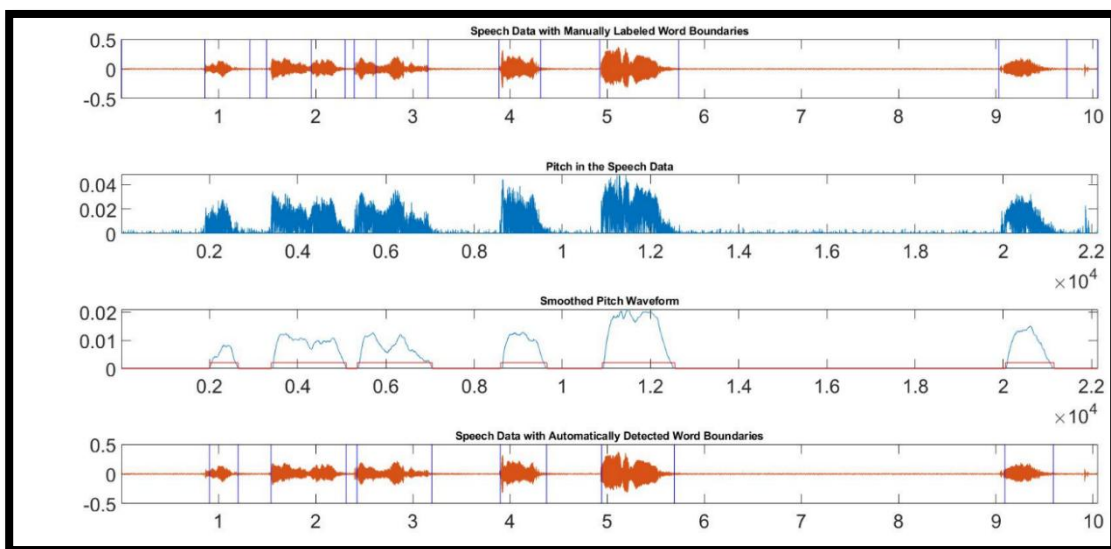


شکل ۴-۸: تشخیص مرز کلمه با استفاده از نرخ عبور از صفر برای داده دوم با طول پنجره ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه با مقایسه نتایج بدست آمده در اشکال (۴-۵) تا (۴-۸)، مشاهده می‌شود که ابعاد پنجره و میزان شیفت زمانی در دقت جداسازی کلمات بسیار تأثیرگذار است. با بررسی ابعاد پنجره در فصل دوم مشخص شد ابعاد پنجره در نتیجه نهایی تأثیر دارد. اگر ابعاد پنجره بدون آن که میزان شیفت زمانی تغییر کند، افزایش یابد، زمان پردازش بیشتر می‌گردد بدون آن که دقت جداسازی تغییر کند. اگر شیفت زمانی خیلی کوچک انتخاب شود، زمان پردازش به شدت افزایش یافته و دقت جداسازی کلمات نیز کاهش می‌یابد بدین صورت که مرزهای بیشتری برای کلمات پیدا می‌کند. از طرف دیگر اگر میزان شیفت زمانی خیلی بزرگ (به نسبت ابعاد پنجره) انتخاب گردد، در اینصورت دقت جداسازی کلمات کاهش می‌یابد و هر چند کلمه را به عنوان یک بخش در نظر می‌گیرد. به همین خاطر انتخاب ابعاد پنجره و مهم تر از آن میزان شیفت زمانی در دقت جداسازی مرز کلمات تأثیر بسیار زیادی دارد. در شکل (۴-۸) برای داده دوم مشاهده می‌شود.

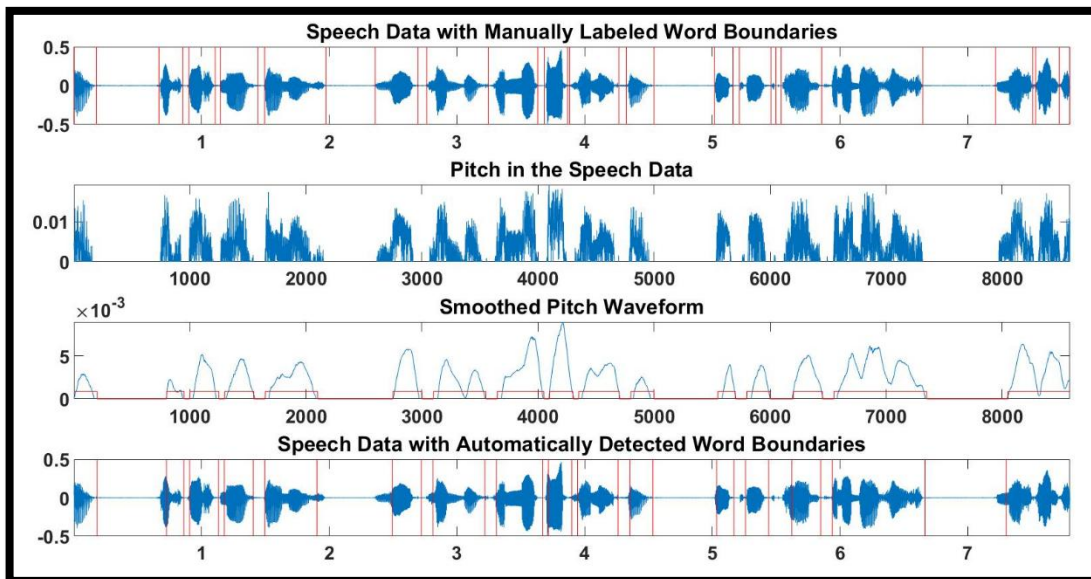
۴-۴-۲- نتیجه شبیه سازی با استفاده از روش گام

همانطور که در فصل دو گفته شد، روش‌های مختلفی برای یافتن پیچ سیگنال گفتار استفاده می‌شود. در این پایان‌نامه برای به دست آوردن پیچ از خودهمبستگی سیگنال استفاده شده است. ابتدا به منظور

بررسی زمان-کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره به طول ۳۲۰ نمونه و با شیفتم زمانی ۲۰ نمونه، فریم بندی می گردد. سپس با استفاده از خودهمبستگی گام فریم‌های واکنار استخراج شده، محاسبه می شود. خود همبستگی سیگنال یعنی محاسبه میزان شباهت سیگنال و مقدار شیفتم یافته‌اش. در مرحله بعد به منظور کم کردن تغییرات خروجی، مقدار DC خروجی این مرحله حذف می شود. برای نرم تر کردن تغییرات منحنی مربوطه، مقادیر بدست آمده توسط فیلتر باتروث فیلتر می گردد و در نهایت با توجه به میزان آستانه‌ای که از روش سعی و خطا به دست آمده، محدوده هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می گردند همانطور که در شکل (۹-۴) و (۱۰-۴) با مقایسه نمودار مربوط به نتایج جداسازی کلمات بصورت دستی و استفاده از الگوریتم فرکانس گام مشاهده می گردد، الگوریتم مبتنی بر فرکانس گام برای هر دو داده توانسته است بیشتر مرزهای کلمات را پیدا کند، اما برخی از مرزها را پیدا نکرده است و در آن قسمت‌ها چند کلمه با هم در یک محدوده قرار گرفته اند.



شکل ۹-۴: تشخیص مرز کلمه با استفاده از روش گام برای داده اول



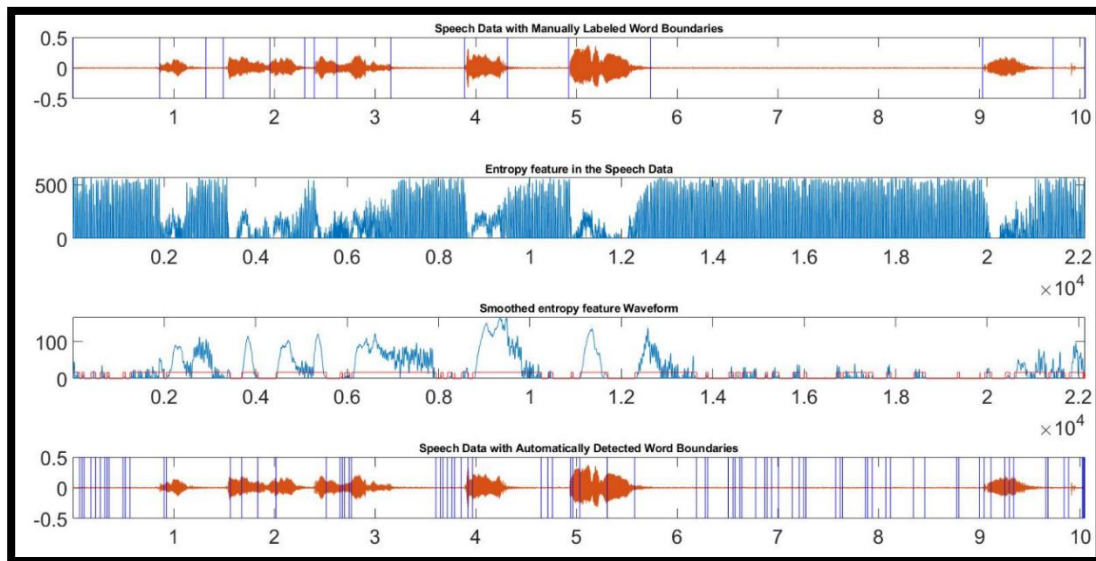
شکل ۴-۱۰: تشخیص مرز کلمه با استفاده از روش گام برای داده دوم

۴-۳-۴- نتیجه شبیه‌سازی روش آنتروپی

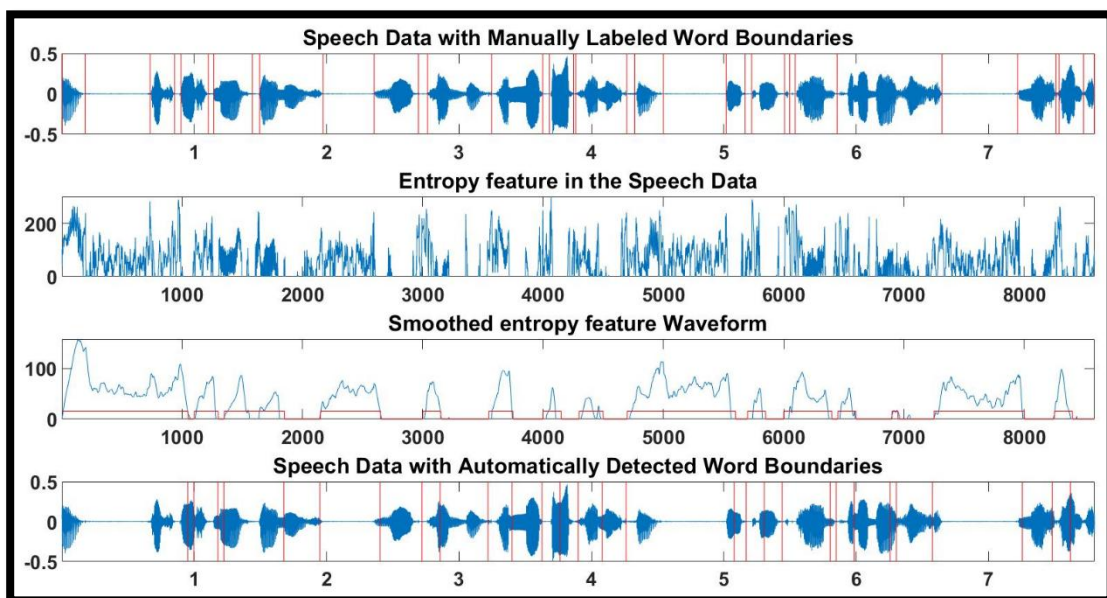
برای تشخیص مرز کلمات با استفاده از آنتروپی هر فریم برای بررسی زمان-کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره به طول ۳۲۰ نمونه و با شیفت زمانی ۲۰ نمونه، فریم بندی می‌گردد. سپس مقادیر فریم‌ها به ماکزیمم مقدار در هر فریم نرمالیزه می‌گردند. در مرحله بعد مقدار آنتروپی مربوط به هر فریم محاسبه شده و سپس مقدار DC خروجی این مرحله حذف می‌شود. در اینجا از سیگنال اصلی برای محاسبه آنتروپی استفاده شده است. به منظور نرم‌تر کردن تغییرات منحنی مربوطه، مقادیر بدست آمده توسط فیلتر باتروث فیلتر می‌گردد و سپس مقدار DC خروجی این مرحله حذف می‌شود و در نهایت با توجه به میزان آستانه تعریف شده، محدوده‌هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می‌گردند.

همانطور که در شکل (۴-۱۱) و (۴-۱۲) با مقایسه نمودار مربوط به نتایج جداسازی کلمات به صورت دستی و استفاده از الگوریتم مربوط به آنتروپی مشاهده می‌گردد، الگوریتم مبتنی بر آنتروپی در مشخص کردن مرز کلمات برای هر دو داده نتوانسته است موفق عمل کند و به اشتباه مرزهای زیادی را برای کلمات پیدا کرده است. به همین جهت این روش نمی‌تواند روش قوی‌ای برای شناسایی مرز کلمات

باشد.



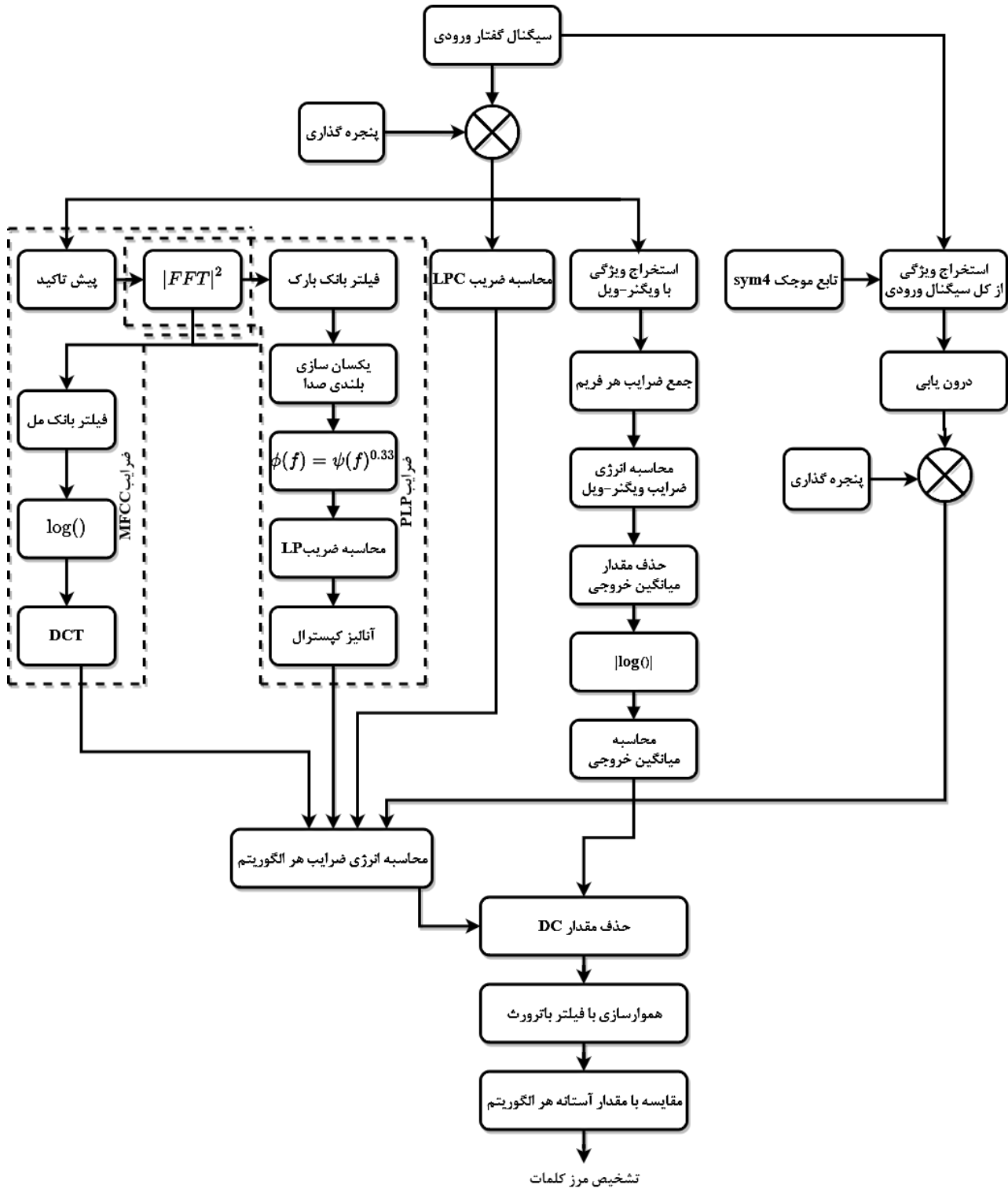
شکل ۴-۱۱: تشخیص مرز کلمه با استفاده از روش آنتروپی برای داده اول



شکل ۴-۱۲: تشخیص مرز کلمه با استفاده از روش آنتروپی برای داده دوم

۴-۵- بلوک دیاگرام کلی روش‌های پیشنهادی

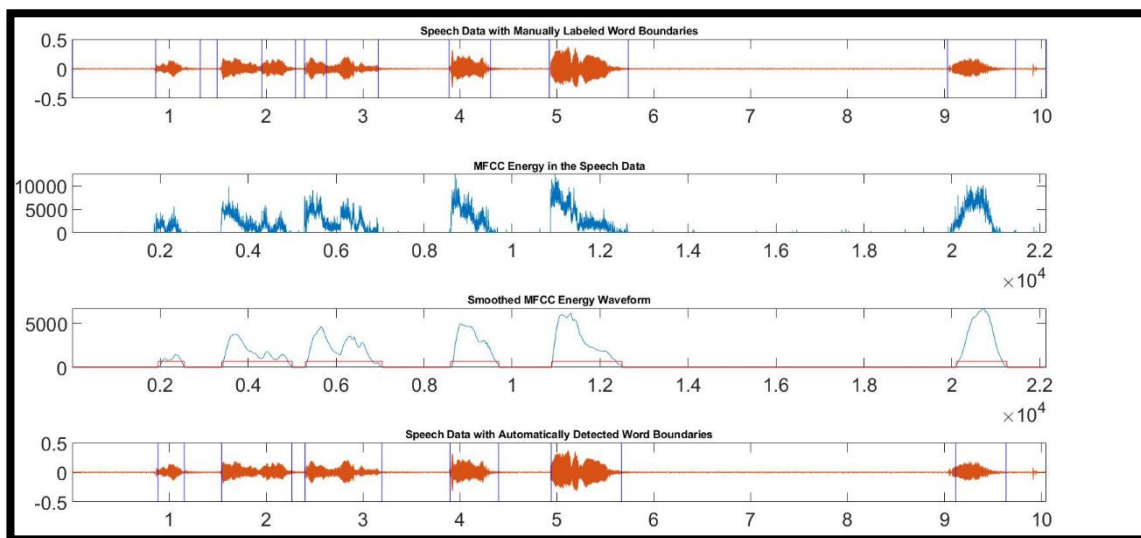
در شکل (۴-۹) بلوک دیاگرام پیشنهادی برای تشخیص مرز کلمات نشان داده شده است.



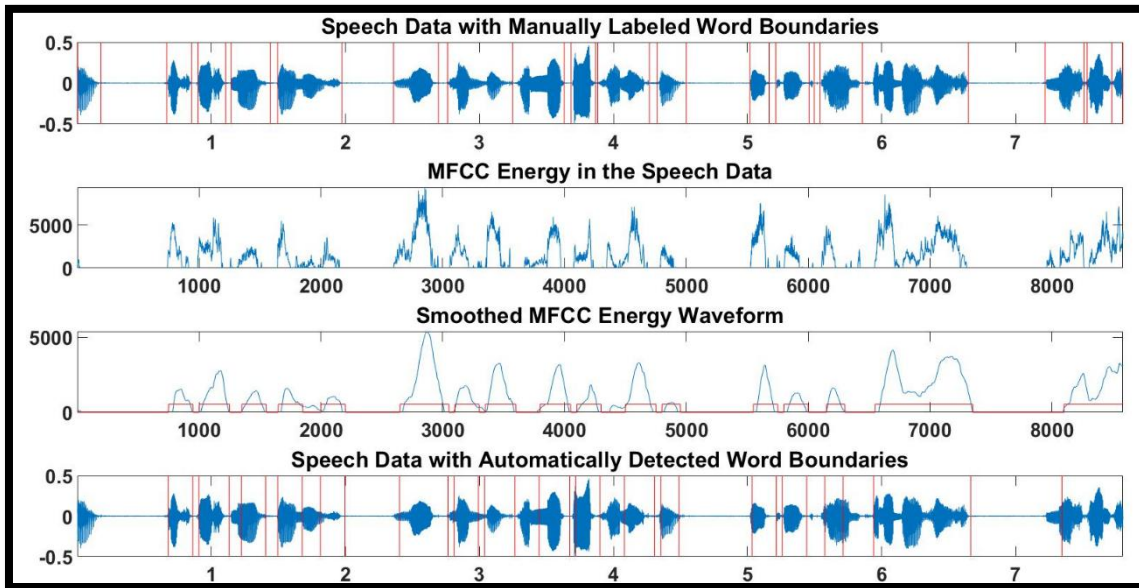
شکل ۴-۱۳: بلوک دیاگرام پیشنهادی برای تشخیص مرز کلمات

۴-۵-۱- نتیجه شبیه‌سازی روش MFCC

در الگوریتم تشخیص مرز کلمات با استفاده از روش استخراج ویژگی MFCC به منظور بررسی زمان- کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره به طول ۳۲۰ نمونه و با شیفت زمانی ۲۰ نمونه، فریم بندی می‌گردد. سپس توسط روش MFCC از هر یک از فریم‌ها ۱۳ ویژگی استخراج می‌گردد و در مرحله بعد انرژی ضرایب MFCC مربوط به هر فریم محاسبه می‌شود. برای محاسبه انرژی ضرایب تک تک ضرایب را به توان دو رسانده باهم جمع می‌شوند. سپس برای کم کردن تغییرات خروجی مقدار DC محاسبه می‌شود و از خروجی سیگنال حذف می‌شود. برای به دست آوردن مقدار DC جمع ضرایب ویژگی‌های انرژی تقسیم بر تعداد ضرایب می‌شود. سپس به منظور نرم‌تر کردن تغییرات منحنی مربوطه، مقادیر به دست آمده توسط فیلتر باتروث فیلتر می‌گردد و در نهایت با توجه به میزان آستانه تعریف شده، محدوده‌هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می‌گردند.



شکل ۴-۱۴: تشخیص مرز کلمه با استفاده از روش MFCC برای داده اول



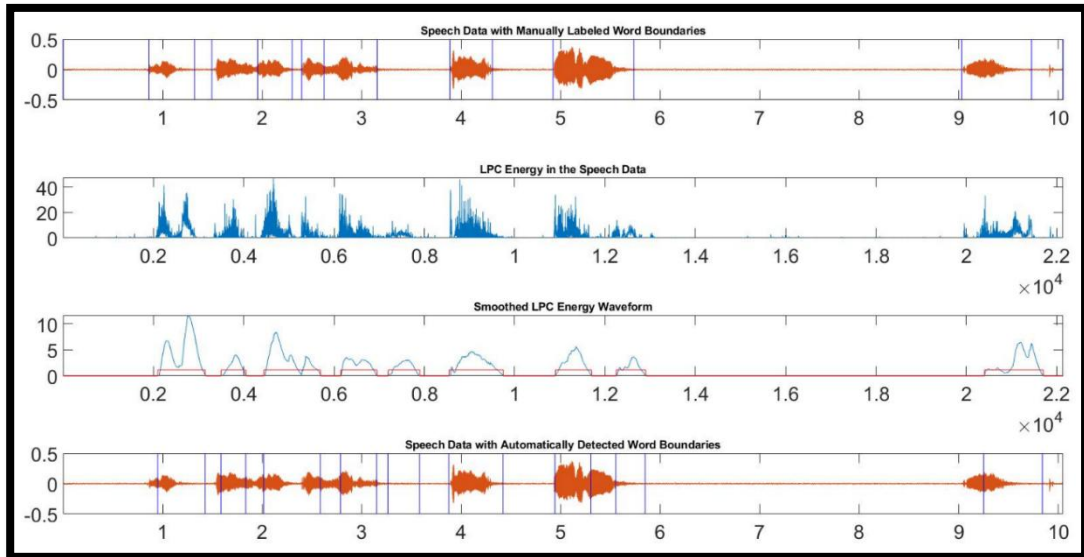
شکل ۴-۱۵: تشخیص مرز کلمه با استفاده از روش MFCC برای داده دوم

همانطور که در شکل (۴-۱۴) و (۴-۱۵) با مقایسه نمودار مربوط به نتایج جداسازی کلمات به صورت دستی و استفاده از الگوریتم استخراج ویژگی MFCC مشاهده می‌گردد، الگوریتم مبتنی بر انرژی ضرایب فرکانس مل توانسته است بیشتر مرزهای کلمات را پیدا کند اما برخی از مرزها را پیدا نکرده است و در آن قسمت‌ها چند کلمه با هم در یک محدوده قرار گرفته‌اند. اما برای داده اول بسیار بهتر توانسته است مرزهای کلمات را تشخیص دهد.

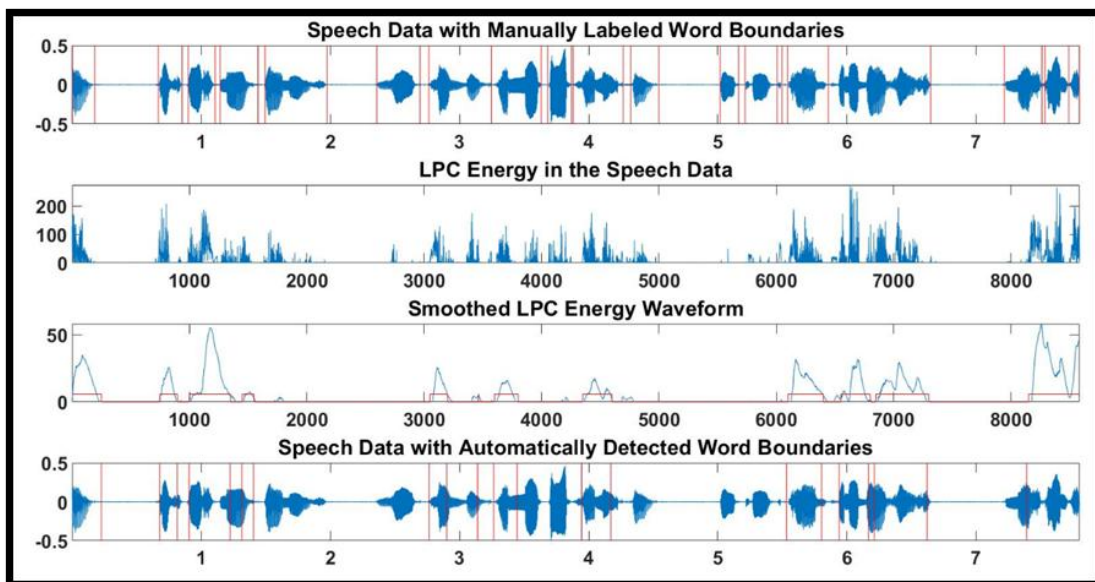
۴-۵-۲- نتیجه شبیه‌سازی روش LPC

الگوریتم جداسازی مرز کلمات با استفاده از روش استخراج ویژگی LPC بدین صورت است که ابتدا به منظور بررسی زمان کوتاه سیگنال-گفتار، سیگنال ورودی با ضرب در پنجره به طول ۳۲۰ نمونه و با شیفت زمانی ۲۰ نمونه، فریم بندی می‌گردد. سپس توسط روش LPC از هر یک از فریم‌ها ۱۳ ویژگی استخراج می‌گردد و در مرحله بعد انرژی ضرایب LPC مربوط به هر فریم محاسبه می‌شود. برای محاسبه انرژی، ضرایب LPC به توان دو می‌رسند و سپس باهم جمع می‌شوند. سپس مقدار DC خروجی این مرحله حذف می‌شود. برای محاسبه مقدار میانگین ضرایب انرژی باهم جمع می‌شوند و تقسیم به تعداد آن‌ها می‌شوند. به منظور نرم‌تر کردن تغییرات منحنی مربوطه، مقادیر بدست آمده توسط فیلتر

باترورت فیلتر می گردد و در نهایت با توجه به میزان آستانه‌ای که به صورت تجربی به دست آمده است، محدوده‌هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می گردند.



شکل ۴-۱۶: تشخیص مرز کلمه با استفاده از روش LPC برای داده اول



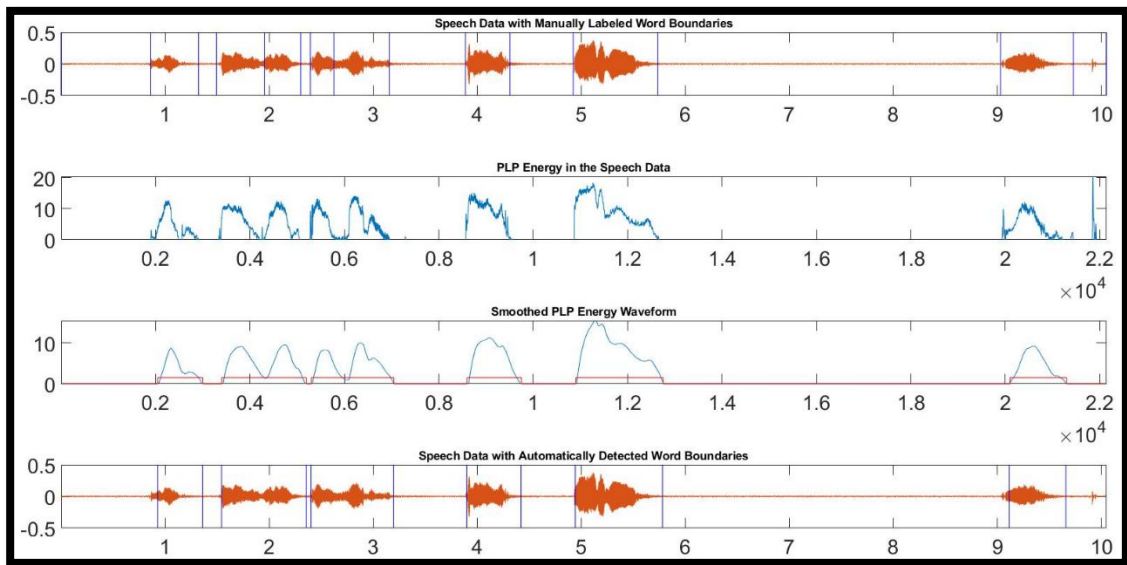
شکل ۴-۱۷: تشخیص مرز کلمه با استفاده از روش LPC برای داده دوم

همانطور که در شکل (۴-۱۶) و (۴-۱۷) با مقایسه نمودار مربوط به نتایج جداسازی کلمات بصورت دستی و استفاده از الگوریتم استخراج ویژگی LPC مشاهده می گردد، الگوریتم مبتنی بر انرژی ضرایب

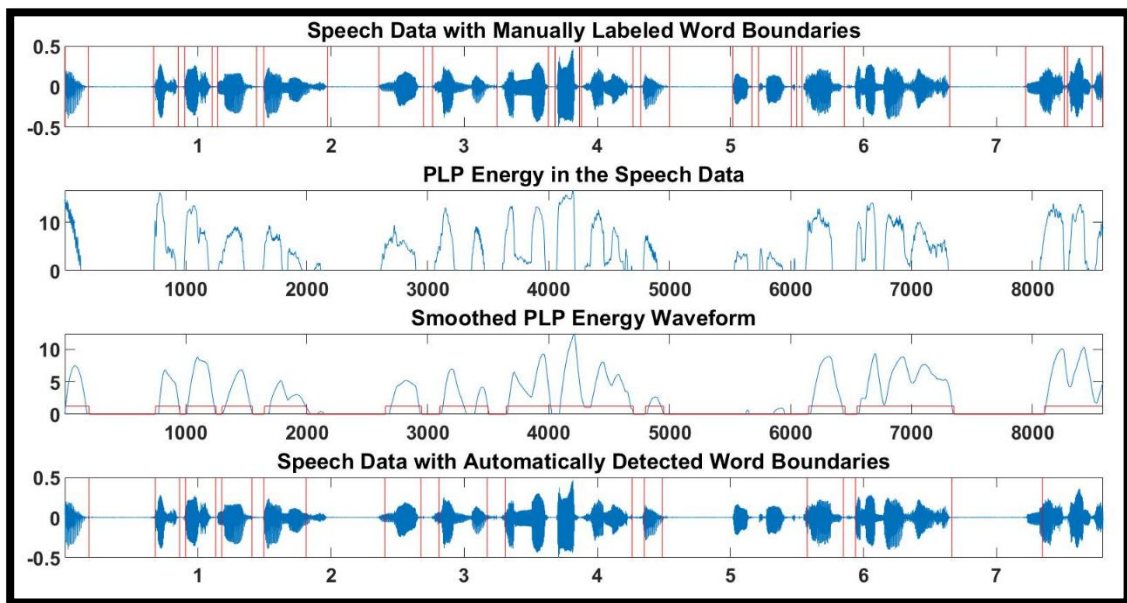
پیشگویی خطی برخی از مرزها را درست و برخی دیگر را به اشتباه شناسایی کرده است. با مقایسه خروجی های LPC و MFCC می توان اینطور بیان نمود که الگوریتم مبتنی بر MFCC توانایی بیشتری در تخمین مرزهای کلمات در قیاس با روش مبتنی بر LPC دارد.

۴-۵-۳- نتیجه شبیه سازی روش PLP

الگوریتم جداسازی مرز کلمات با استفاده از روش استخراج ویژگی PLP بدین صورت است که ابتدا به منظور بررسی زمان-کوتاه سیگنال گفتار، سیگنال ورودی با ضرب در پنجره به طول ۳۲۰ نمونه و با شیفت زمانی ۲۰ نمونه، فریم بندی می گردد. سپس توسط روش PLP از هر یک از فریم ها ۱۳ ویژگی استخراج می گردد و در مرحله بعد انرژی ضرایب PLP مربوط به هر فریم محاسبه می شود. برای محاسبه انرژی هر فریم ویژگی های PLP به دست آمده به توان دو می رسند و با هم جمع می شوند. سپس مقدار DC خروجی این مرحله حذف می شود. مانند قبل برای به دست آوردن مقدار DC مقدار میانگین ضرایب خروجی محاسبه می شود. به منظور نرم تر کردن تغییرات منحنی مربوطه، مقادیر بدست آمده توسط فیلتر باتروث فیلتر می گردد و در نهایت با توجه به میزان آستانه تعریف شده، محدوده هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می گردند. همانطور که در شکل (۴-۱۸) و (۴-۱۹) با مقایسه نمودار مربوط به نتایج جداسازی کلمات بصورت دستی و استفاده از الگوریتم استخراج ویژگی PLP مشاهده می گردد، الگوریتم مبتنی بر انرژی PLP برخی از مرزها را درست و برخی دیگر را به اشتباه شناسایی کرده است. با مقایسه خروجی های مربوط به PLP و LPC می توان اینطور بیان نمود که الگوریتم مبتنی بر PLP توانایی بیشتری در تخمین مرزهای کلمات در قیاس با روش مبتنی بر LPC برای هر دو داده دارد.



شکل ۴-۱۸: تشخیص مرز کلمه با استفاده از روش PLP برای داده اول



شکل ۴-۱۹: تشخیص مرز کلمه با استفاده از روش PLP برای داده دوم

۴-۵-۴ - تئوری کار روش موجک

همانطور که گفته شد سیگنال‌ها تا این قسمت یا در حوزه زمان یا در حوزه فرکانس بودند. برای تحلیل بهتر سیگنال‌ها، در حوزه زمان-فرکانس بررسی می‌شوند. روش موجک و روش ویگنر از روش‌های زمان-فرکانس هستند که نتایج آن‌ها در این قسمت و قسمت بعد آورده شده است.

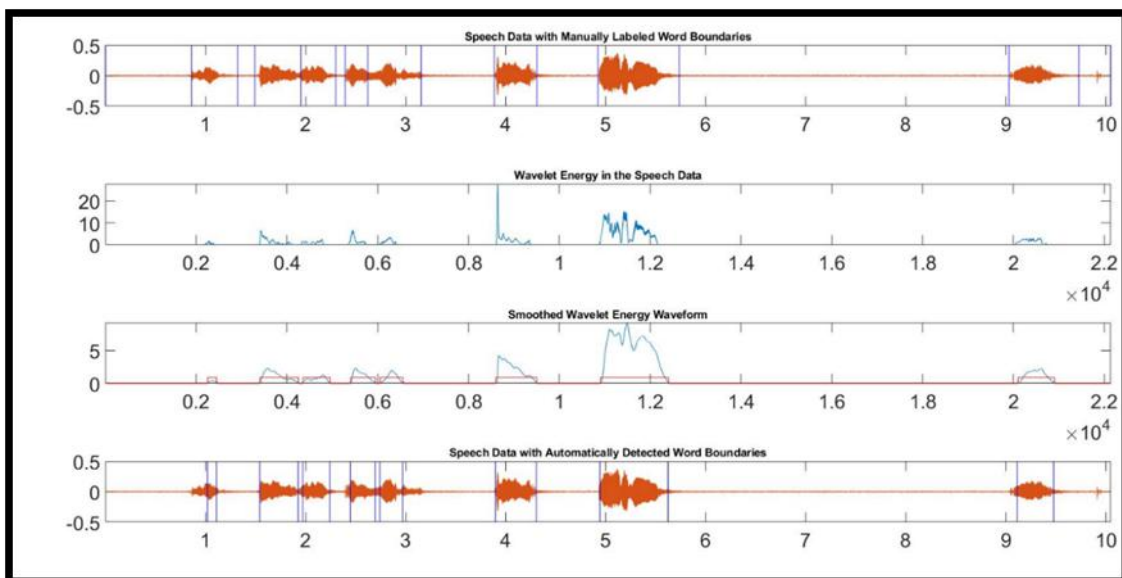
مراحل تشخیص مرز کلمات با استفاده از روش موجک به شرح زیر است:

- ابتدا با استفاده از تبدیل موجک گسسته از کل سیگنال ورودی ویژگی استخراج می‌شود. تابع موجک مورد استفاده sym4 است.
- پس از استخراج ویژگی به این دلیل که فقط از مؤلفه‌های فرکانس پایین یعنی مؤلفه تقریب استفاده می‌شود طول سیگنال نصف شده و نمی‌تواند با روش دستی مقایسه شود. به همین دلیل درون‌یابی صورت می‌گیرد تا طول این سیگنال با سیگنال اصلی برابر شود. برای درون‌یابی دو نمونه قبل و بعد را در نظر گرفته و تقسیم به دو می‌شود.
- در مرحله بعد برای فریم‌بندی، پنجره‌هایی با طول ۳۲۰ نمونه در سیگنال ضرب می‌شود. انرژی هر پنجره که انرژی ویژگی‌های موجک است، محاسبه می‌شود. می‌توان به جای انرژی از معیارهای دیگری مانند آنتروپی نیز استفاده کرد.
- اگر خروجی تغییرات زیادی داشته باشد برای کم کردن تغییرات و پردازش سریع‌تر مقدار DC خروجی را حذف می‌شود. برای این کار میانگین انرژی ویژگی‌ها را محاسبه کرده و از خروجی ویژگی‌ها کم می‌شود تا خروجی تغییرات کمتری داشته باشد.
- سپس برای نرم‌تر کردن تغییرات منحنی و کاهش اثر نویز توسط فیلتر باتروث، فیلتر می‌گردد.
- با استفاده از روش سعی و خطا و جای‌گذاری مقادیر مختلف برای آستانه، بهترین مقدار آستانه برای تشخیص مرز کلمات استفاده می‌شود.
- در مرحله آخر مرزهای کلمات به دست آمده با استفاده از روش موجک با مرزهایی که به صورت دستی به دست آمده‌اند، مقایسه می‌شوند و مقدار تشخیص صحیح مشخص می‌شود.

۴-۵-۴-۱- نتایج شبیه‌سازی روش موجک

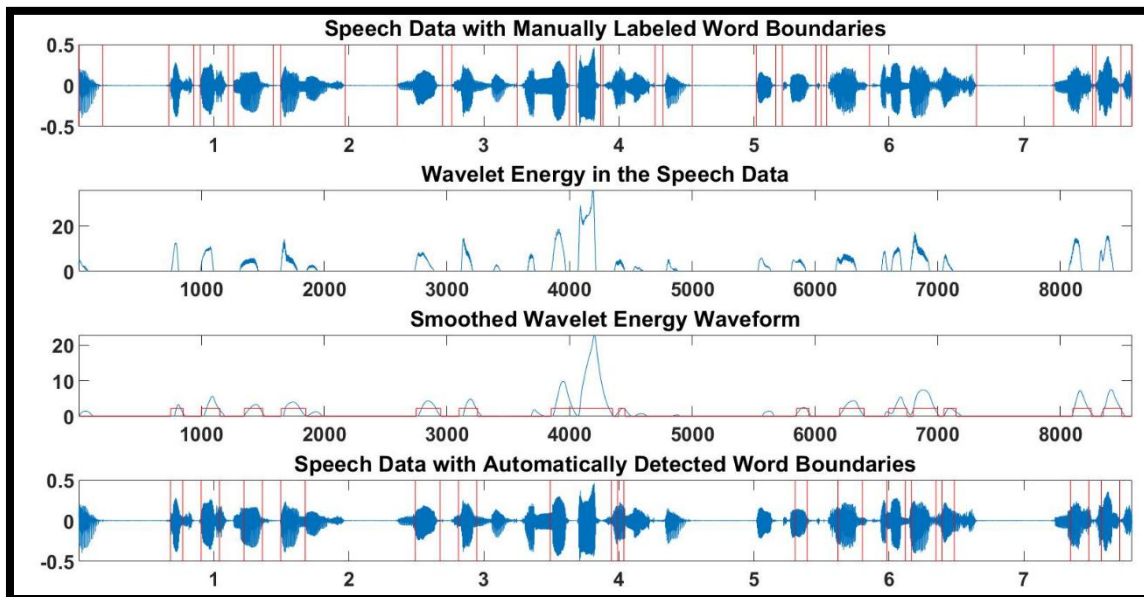
در قسمت قبل الگوریتم جداسازی مرز کلمات با استفاده از روش موجک توضیح داده شد. در این قسمت نتایج شبیه‌سازی با استفاده از تبدیل موجک برای تشخیص مرز کلمات نمایش داده می‌شود. در قسمت اول شکل (۴-۲۰) مرزهای کلمات به صورت دستی مشخص شده‌اند. سپس استخراج ویژگی

انجام شده و درون‌یابی صورت می‌گیرد. در مرحله بعد با فریم بندی خروجی مرحله قبل، انرژی زمان - کوتاه برای هر فریم محاسبه شده و سپس مقدار DC خروجی این مرحله حذف می‌شود و سیگنال توسط فیلتر باترورث، فیلتر می‌گردد. با توجه به میزان آستانه تعریف شده برای داده اول، محدوده هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می‌گردند. در قسمت سوم شکل‌ها بخش‌هایی که بیشتر از مقدار آستانه است با خطوط قرمزی که ۰/۱ ماکزیمم ارتفاع سیگنال است نشان داده شده است. این خطوط همان مرزهای کلمات است که در قسمت چهارم بر روی سیگنال اصلی نشان داده شده است. سپس دو قسمت اول و آخر مقایسه می‌شود و تشخیص‌های صحیح مشخص می‌گردد.



شکل ۴-۲: تشخیص مرز کلمات با استفاده از روش موجک برای داده اول

همانطور که در شکل (۴-۲۰) و (۴-۲۱) با مقایسه نمودار مربوط به نتایج جداسازی کلمات بصورت دستی و استفاده از الگوریتم مبتنی بر موجک مشاهده می‌گردد، الگوریتم مبتنی بر انرژی موجک برخی برای داده‌های اول بسیار بهتر از داده‌های دوم مرزها را مشخص کرده است.



شکل ۴-۲۱: تشخیص مرز کلمات با استفاده از روش موجک برای داده دوم

۴-۵-۵- تئوری کار روش ویگنر-ویل

مراحل تشخیص مرز کلمات با استفاده از روش ویگنر-ویل به شرح زیر است:

- ابتدا سیگنال ورودی با ضرب در پنجره‌ای به طول ۳۲۰ نمونه و شیفت زمانی ۲۰ نمونه فریم‌بندی و پنجره‌گذاری می‌شود.
- سپس تبدیل ویگنر-ویل همه‌ی فریم‌ها محاسبه می‌شود. برای هر فریم تمام توزیع‌ها باهم جمع می‌شود. یعنی برای هر فریم یک بردار وجود دارد.
- در مرحله بعد انرژی ضرایب ویگنر-ویل محاسبه می‌شود. برای محاسبه انرژی، هر ضریب انرژی به توان دو رسیده سپس تمام ضرایب جمع می‌شوند.
- برای کم کردن تغییرات خروجی، میانگین انرژی ویگنر-ویل محاسبه شده و از مقدار خروجی کم می‌شود. برای محاسبه مقدار میانگین تمام ضرایب انرژی با هم جمع می‌شوند و تقسیم بر تعداد ضرایب می‌شود.
- مشکلی که روش ویگنر-ویل دارد این است که اختلاف بین اعداد در سیگنال زیاد است. برای مثال یکی صفر و دیگری ده‌هزار است. در صورتی که تغییرات داده زیاد باشد، برای

دیده شدن تمام اعداد باید مقیاس لگاریتمی شود. بعد از استخراج لگاریتم چون اعداد مثبت و منفی می‌شوند، برای جلوگیری از این مشکل قدر مطلق اعداد نیز محاسبه می‌شود.

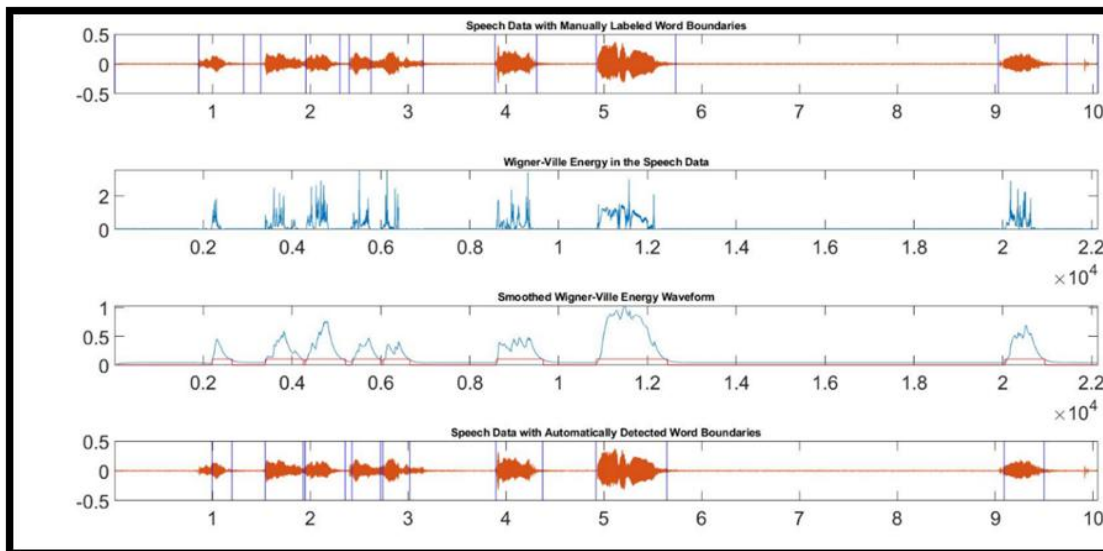
- در مرحله بعد میانگین را از خروجی داده‌های لگاریتم گرفته شده کم می‌شود. به دلیل این که ممکن است بعضی از داده‌های لگاریتم گرفته شده از مقدار میانگین کمتر شود بار دیگر قدر مطلق گرفته می‌شود تا نتیجه همیشه مثبت باشد.
- در مرحله بعد برای تشخیص بهتر مرزهای کلمه هموارسازی انجام می‌شود.
- با توجه به مقدار آستانه محدوده‌ای از سیگنال که بیشتر از مقدار آستانه باشد به عنوان مرز کلمه در نظر گرفته می‌شود.

۴-۵-۵-۱- نتیجه شبیه‌سازی روش ویگنر-ویل

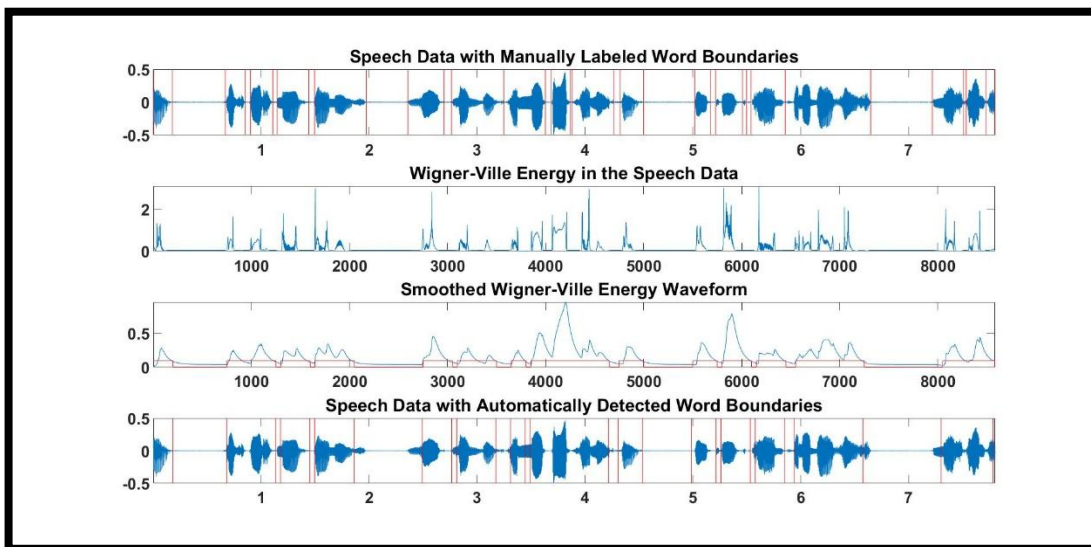
در قسمت قبل مراحل تشخیص مرز کلمات با استفاده از روش ویگنر-ویل توضیح داده شد. در این قسمت نتایج شبیه‌سازی با استفاده از توزیع ویگنر-ویل برای تشخیص مرز کلمات نمایش داده می‌شود. در قسمت اول شکل (۴-۲۲) مرزهای کلمات به صورت دستی مشخص شده‌اند. ابتدا سیگنال فریم‌بندی شده و سپس با استفاده از روش ویگنر-ویل استخراج ویژگی انجام می‌شود. مقدار DC خروجی این مرحله حذف می‌شود و سیگنال توسط فیلتر باترورث، فیلتر می‌گردد. با توجه به میزان آستانه تعریف شده، محدوده‌هایی از منحنی که دارای مقدار بزرگتری از سطح آستانه است به عنوان مرز کلمات انتخاب می‌گردند. در قسمت سوم شکل (۴-۲۲) بخش‌هایی که بیشتر از مقدار آستانه است با خطوط قرمزی که ۰/۱ ماکزیمم ارتفاع سیگنال است نشان داده شده است. این خطوط همان مرزهای کلمات است که در قسمت چهارم بر روی سیگنال اصلی نشان داده شده است. سپس دو قسمت اول و آخر مقایسه می‌شود و تشخیص‌های صحیح مشخص می‌گردد.

همانطور که در شکل (۴-۲۲) و (۴-۲۳) با مقایسه نمودار مربوط به نتایج جداسازی کلمات بصورت دستی و استفاده از الگوریتم مبتنی بر لگاریتم انرژی ضرایب ویگنر-ویل مشاهده می‌گردد که الگوریتم پیشنهادی در قیاس با سایر روش‌ها برای هر دو داده توانسته است تعداد بیشتری مرز را به درستی

تخمین بزند. با توجه به نتایج بدست آمده می توان اینطور بیان نمود که الگوریتم پیشنهادی مناسبترین و قوی ترین روش از بین روش های بررسی شده برای شناسایی مرز کلمات است.



شکل ۴-۲۲: تشخیص مرز کلمات با استفاده از روش ویگنر-ویل برای داده اول

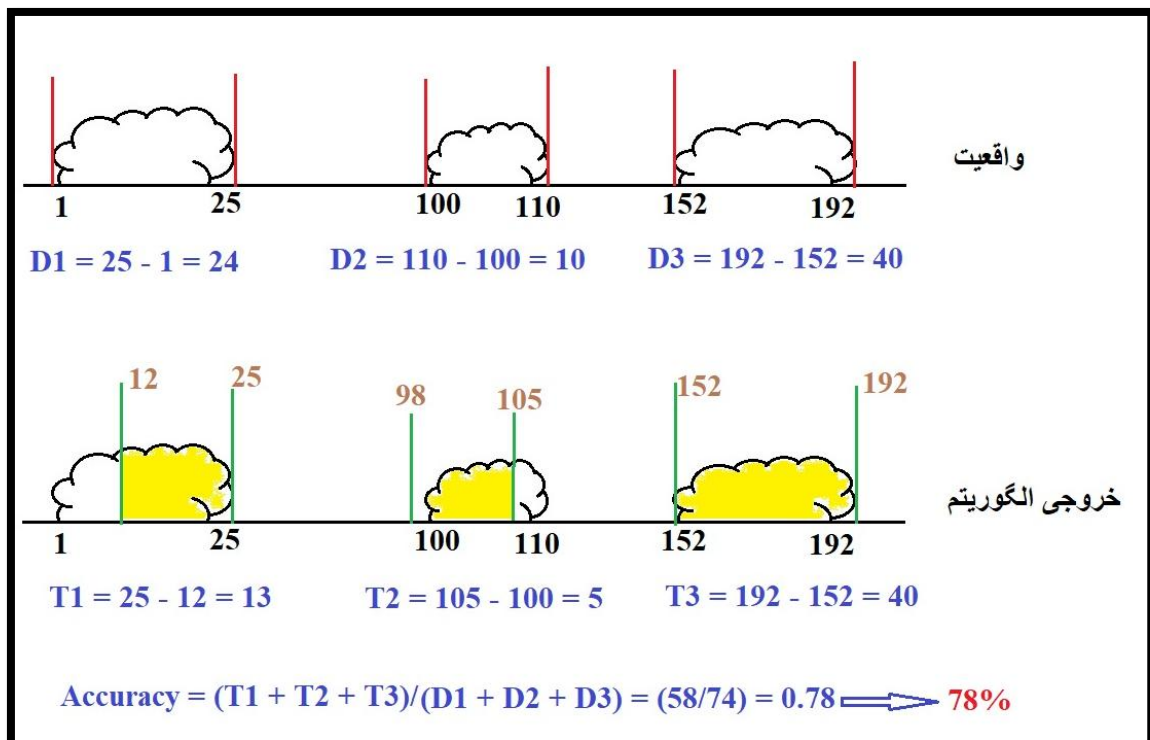


شکل ۴-۲۳: تشخیص مرز کلمات با استفاده از روش ویگنر-ویل برای داده دوم

۴-۶- نتیجه گیری

در بخش های قبل الگوریتم تشخیص مرز کلمات توضیح داده شده و نتایج شبیه سازی روش های مختلف نشان داده شده است. برای به دست آوردن دقت تشخیص صحیح در هر یک از روش ها مانند

شکل (۴-۲۵) عمل می‌کنیم. این روش، شهودی است و دارای دقت خوبی می‌باشد. مرزهایی که با استفاده از اطلاعات موجود در پایگاه داده مشخص می‌شوند در قسمت واقعیت نشان داده شده است. با استفاده از روش‌های مختلف مرزهای کلمه به دست می‌آیند و در قسمت خروجی الگوریتم نشان داده می‌شود. برای مثال در شکل (۴-۲۴) ۳ مرز مشخص شده است. مرز قسمت اول از نمونه‌ی ۱ تا نمونه ۲۵ است. اما خروجی الگوریتم ابتدای مرز را نمونه ۱۲ و انتهای آن را نمونه ۲۵ تشخیص داده است. یعنی الگوریتم تا حدودی مرز را درست در تشخیص داده است. با توجه به توضیحات داخل شکل دقت تشخیص، مشخص می‌شود. پس در این روش مرز تشخیص داده کلاً اشتباه در نظر گرفته نمی‌شود و درصدی از آن صحیح است. در صورتی که بازه‌ای تشخیص داده شود که در محدود صحیح نباشد در محاسبات در نظر گرفته نمی‌شود.



شکل ۴-۲۴: روش به دست آوردن دقت تشخیص مرز کلمات

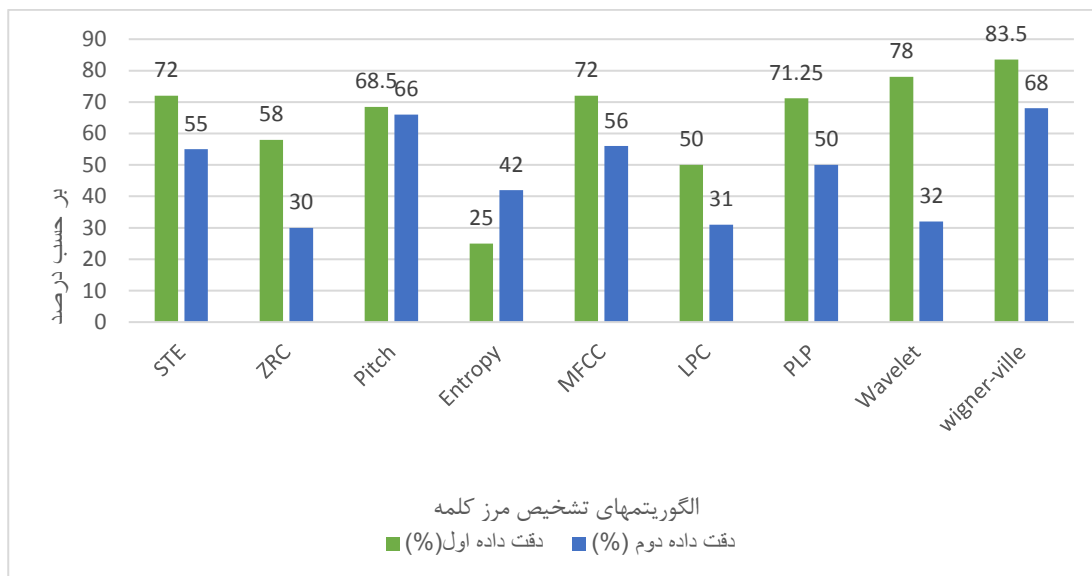
بعد از به دست آوردن دقت تشخیص هر روش مشاهده می‌شود که روش ویگنر-ویل که روش استخراج ویژگی از نوع زمان-فرکانس است بهترین روشی است که برای تشخیص مرز کلمات در سیگنال‌های گفتاری استفاده می‌شود. پس تحلیل سیگنال در هر دو حوزه زمان و فرکانس باعث

تشخیص بهتر مرز کلمات می‌شود. همچنین روش موجک که روش دیگری از روش‌های زمان-فرکانس است دقت بالایی در تشخیص مرز کلمات در داده‌های با پیچیدگی کمتر دارد. در نتیجه سیگنال گفتار که یک سیگنال ناپایستا است و مشخصات آماری آن با زمان تغییر می‌کند، تحلیل این سیگنال در هر دو حوزه زمان و فرکانس باعث تشخیص بهتر مرز کلمات می‌شود. آنالیزهای زمانی در تشخیص مرزهای کلمات خیلی قوی نبودند. برای مثال در روش آنتروپی از سیگنال خام استفاده شده است که نتیجه مطلوبی نداشته است. اگر ابتدا از سیگنال تبدیل فوریه گرفته شود و سپس با استفاده از روش آنتروپی مرز کلمات مشخص شود، دقت نتیجه بهتر می‌شد. در جدول (۴-۴) دقت تشخیص هر یک از روش‌ها نشان داده شده است.

جدول ۴-۳: نتایج روش‌های تشخیص مرز کلمه

روش	STE	ZRC	Pitch	Entropy	MFCC	LPC	PLP	Wavelet	Wigner-ville
دقت برای داده اول (%)	۷۲	۵۸	۶۸/۵	۲۵	۷۲	۵۰	۷۱/۲۵	۷۸	۸۳/۵
دقت برای داده دوم (%)	۵۵	۳۰	۵۹	۴۲	۵۶	۳۱	۳۲	۴۲	۶۸

طبق نتایج ارائه شده، در داده‌هایی با پیچیدگی کمتر، الگوریتم‌ها عملکرد تقریباً مشابهی دارند. در واقع میزان مکث‌های بین کلمات و نحوه تلفظ گوینده در دقت تشخیص مرزها مؤثر است. با این حال، هنگامی که از داده‌هایی با پیچیدگی بیشتر استفاده می‌شود، شاهد افت کارایی نسبتاً زیاد روش‌های قبلی در مقایسه با روش ویگنر-ویل هستیم.



شکل ۴-۲۵: نمودارهای دقت برای داده اول و داده دوم

فصل پنجم:

نتیجه‌گیری و

پیشنهادات

۵-۱- جمع بندی

در این تحقیق، هدف تشخیص مرزهای کلمه با استفاده از روش‌های مختلف استخراج ویژگی است. تشخیص مرز کلمه یعنی مشخص کردن ابتدا و انتها یک کلمه در یک سیگنال گفتار. تشخیص مرزهای کلمات در بسیاری از کاربرها مانند بازشناسی گفتار، سنتز گفتار، تبدیل متن به گفتار و... کاربرد دارد. از دونوع داده در این پایان‌نامه استفاده شده است. پایگاه داده TIMIT و داده‌ای که کلمات را به طور واضح بیان می‌کند. الگوریتم‌های استخراج ویژگی به سه دسته آنالیز زمانی، آنالیز طیفی و آنالیز زمان-فرکانس تقسیم می‌شوند. در فصل دوم روش‌های مرسوم که برای تشخیص مرز کلمه استفاده شده است و شامل انرژی، نرخ عبور از صفر، گام و آنتروپی، بررسی شدند. ابتدا مراحل پیش پردازش سیگنال گفتار به دلیل تغییر مشخصه‌های سیگنال گفتار در طول زمان و به عبارتی ناپیدا بودن آن صورت می‌گیرد. در واقع استخراج ویژگی از یک حوزه زمانی نسبتاً بزرگ، اطلاعات معتبری را در اختیار نمی‌گذارد. به همین منظور سیگنال گفتار فریم‌بندی و پنجره‌گذاری می‌شود. پنجره مورد نظر همینگ و طول پنجره ۳۲۰ نمونه و مقدار شیفیت پنجره ۲۰ نمونه می‌باشد. پنجره بسیار مهم است زیرا که حاشیه‌های یک فریم در کم یا زیاد شدن خطا مؤثرند. با استفاده از الگوریتم مورد نظر ویژگی‌ها استخراج می‌شوند. سپس توسط فیلتر باترورث که یک فیلتر پایین‌گذر و درجه یک است به منظور کم کردن نویز و افزایش کارایی سیستم، هموارسازی سیگنال انجام می‌شود. سپس برای تشخیص مرزها آستانه‌ای در نظر گرفته می‌شود. روش به دست آوردن آستانه به صورت تجربی و سعی و خطا است. اگر مقدار آستانه بالا در نظر گرفته شود باعث از دست رفتن تعداد زیادی از مرزهای کلمات می‌شود. پس چند مقدار آستانه آزمایش می‌شود و بهترین آن که تعداد بیشتر از مرزها را تشخیص انتخاب می‌شود. در روش مورد نظر، اگر ویژگی استخراج شده مقدارش از آستانه بیشتر شود به عنوان مرز کلمات در نظر گرفته می‌شود. سپس مرزهای تشخیص داده شده با الگوریتم مورد نظر با مرزهایی که به صورت دستی و از اطلاعات موجود در پایگاه داده مشخص شده‌اند، مقایسه می‌شود و میزان تشخیص صحیح مرزهای

کلمات مشخص می‌شود. در این روش‌ها به دلیل وجود نویز در محیط، استفاده از برخی از این معیارها چندان مناسب نمی‌باشد. همچنین عدم وجود مکث کافی باعث خطا در شناسایی مرزها می‌شود. گوینده باید کلمات را به وضوح تلفظ کند.

در فصل سوم از روش‌های طیفی و روش‌های زمان-فرکانس برای شناسایی مرزهای کلمات استفاده شده است. MFCC، LPC، PLP، موجک و ویگنر-ویل از جمله‌ی این روش‌هاست. در سه روش اول پس از فریم بندی و پنجره‌گذاری از هر فریم ۱۳ ویژگی استخراج می‌شود و میزان انرژی ضرایب هر الگوریتم محاسبه می‌شود. در الگوریتم موجک با استفاده از تابع موجک مؤلفه تقریب به دست می‌آید و سپس به منظور همسان‌سازی طول مؤلفه تقریب با سیگنال اصلی درون یابی صورت می‌گیرد و انرژی ضرایب موجک به دست می‌آید. در ویگنر-ویل نیز ضرایب ویگنر-ویل از هر فریم استخراج می‌شود و ضرایب در هر فریم با هم جمع می‌شوند و مقدار انرژی ضرایب محاسبه می‌شود. پس از محاسبه ضرایب انرژی و حذف مقدار DC سیگنال، مقدار آستانه در هر الگوریتم اعمال می‌شود. در این الگوریتم‌ها نیز مقدار آستانه با سعی و خطا و روش تجربی به دست آمده که ضریبی از ماکزیمم ارتفاع سیگنال هموار شده می‌باشد. مرزها مشخص می‌شود و با مقدار مرزهایی که به روش دستی مشخص شده‌اند مقایسه می‌گردند. طبق نتایج به دست آمده روش ویگنر-ویل بیشترین مقدار تشخیص صحیح را برای هر دو داده دارد. در نهایت نتایج شبیه‌سازی در فصل چهارم نشان داده شده است.

۵-۲- پیشنهادات برای ادامه کار

در این قسمت پیشنهاداتی برای تشخیص مرز کلمه و بهبود نتایج داده شده است:

- یکی از چالش‌های اصلی تشخیص مرز کلمات، انتخاب مقدار آستانه می‌باشد. اگر بتوان برای تمام روش‌ها آستانه را با روش مشخصی و به صورت یکسان محاسبه کرد، بهتر می‌توان الگوریتم‌ها را با یکدیگر مقایسه کرد. حتی می‌توان دو مقدار برای آستانه در نظر گرفت. یکی برای شناسایی ابتدای کلمات و دیگر برای شناسایی انتهای کلمات.

- همچنین برای هموار سازی سیگنال‌های خروجی می‌توان از فیلترهای دیگر با درجه‌های متنوع استفاده کرد تا تأثیر آن‌ها در نتیجه نهایی مشخص شود.
- می‌توان از الگوریتم‌های زمان-فرکانس دیگر همچون روش تبدیل فوریه کسری و یا الگوریتم‌های یادگیری همچون شبکه‌های عصبی، فازی و احتمالاتی آزمایش کرد.
- همانطور که در فصل دو مشاهده شده ترکیب روش ZRC با سایر روش‌ها باعث بهبود نتایج شده بود. می‌توان الگوریتم‌های مختلف مانند MFCC، ویگنر-ویل و آنتروپی را باهم ترکیب کرد تا باعث بهبود دقت شود.

منابع

- [1] A. Caranica, H. Cucu, A. Buzo and C. Burileanu, "On the design of an automatic speech recognition system for Romanian language," *Journal of Control Engineering and Applied Informatics*, vol. 18, no. 2, pp. 65-76, 2016.
- [۲] آیت، مبانی پردازش سیگنال گفتار، انتشارات دانشگاه پیام نور، ۱۳۸۷.
- [3] L. Lamel, L. Rabiner and A. Rosenberg, "An Improved Endpoint Detector for Isolated Word Recognition," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 29, no. 4, pp. 777-785, 1981.
- [4] J.-C. Junqua, B. Mak and B. Reaves, "A Robust Algorithm for Word Boundary Detection in the Presence of Noise," *IEEE Transactions on speech and audio processing*, vol. 2, no. 3, pp. 406-412, 1994.
- [5] A. D. Vijayendra and V. K. Thakar, "Word boundary detection for Gujarati speech recognition using in-ear microphone," in *2016 1st India International Conference on Information Processing (IICIP)*, Delhi, India, 2016.
- [6] A. Agarwal, A. Jain and N. Prakash, "Word Boundary Detection in Continuous Speech based on Suprasegmental Features for Hindi Language," in *2010 2nd International Conference on Signal Processing Systems*, Dalian, China, 2010.
- [7] Z. Hu, J. Schalkwyk, E. Barnard and R. Cole, "Speech recognition using syllable-like units," in *Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP '96*, Philadelphia, PA, USA, 1996.
- [8] D. Enqing, L. Guizhong, Z. Yatong and C. Yu, "Voice activity detection based on short-time energy and noise spectrum adaptation," in *6th International Conference on Signal Processing, 2002*, Beijing, China, 2002.
- [9] D. Arifianto, "Dual parameters for voiced-unvoiced speech signal determination," in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing - ICASSP '07*, Honolulu, HI, USA, 2007.
- [10] C. Guanghua, L. Junhai and Y. Jun, "An Improved Method of Endpoints Detection Based on Energy-Frequency-Value," in *Conference on High Density Microsystem Design and Packaging and Component Failure Analysis, 2006. HDP'06*, Shanghai, China, 2006.
- [11] M. Jalil, F. A. Butt and A. Malik, "Short-time energy, magnitude, zero crossing rate and autocorrelation measurement for discriminating voiced and unvoiced segments of speech signals," in *2013 The International Conference on Technological*

Advances in Electrical, Electronics and Computer Engineering (TAECE), Konya, Turkey, 2013.

- [12] S. Kumar, S. Phadikar and K. Majumder, "Modified segmentation algorithm based on short term energy & Zero crossing rate for maithili speech signal," in *2016 International Conference on Accessibility to Digital World (ICADW)*, Guwahati, India, 2016.
- [13] T. H. Zaw and N. War, "The combination of spectral entropy, zero crossing rate, short time energy and linear prediction error for voice activity detection," in *2017 20th International Conference of Computer and Information Technology (ICCIT)*, Dhaka, Bangladesh, 2017.
- [14] L. Hui, B. Dai and L. Wei, "A pitch detection algorithm based on amdf and acf," in *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, Toulouse, France, 2006.
- [15] G. R. Rao and J. Srichland, "Word boundary detection using pitch variations," in *Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP '96*, Philadelphia, PA, USA, 1996.
- [16] J. Shen, J. Hung and L. Lee, "Robust Entropy-based Endpoint Detection for Speech Recognition in Noisy Environments," *ICSLP*, vol. 98, pp. 232-235, 1998.
- [17] L. Jin and J. Cheng, "An improved speech endpoint detection based on spectral subtraction and adaptive sub-band spectral entropy," in *2010 International Conference on Intelligent Computation Technology and Automation*, Changsha, China, 2010 .
- [18] B.-F. Wu and K.-C. Wang, "Robust endpoint detection algorithm based on the adaptive band-partitioning spectral entropy in adverse environments," *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 5, pp. 762-775, 2005.
- [19] Y. Zhang, K. Wang and B. Yan, "Speech endpoint detection algorithm with low signal-to-noise based on improved conventional spectral entropy," in *2016 12th World Congress on Intelligent Control and Automation (WCICA)*, Guilin, China, 2016.
- [20] P.-S. Antonia, *Applications in Time-Frequency Signal Processing*, CRC press, 2018.
- [21] R. D. Patterson, "Auditory filter shape," *The Journal of the Acoustical Society of America*, vol. 55, no. 4, pp. 802-809, 1974.
- [22] M.-A. Nasr, M. Abd-Elnaby, A.-S. El-Fishawy, S. El-Rabaie and F.-E. A. El-Samie, "Speaker identification based on normalized pitch frequency and Mel Frequency Cepstral Coefficients," *International Journal of Speech Technology*, vol. 21, pp. 941-951, 2018.

- [23] I. Daly, Z. Hajaiej and A. Gharsallah, "Speech analysis in search of speakers with MFCC , PLP , Jitter and Shimmer," in *2017 International Conference on Advanced Systems and Electric Technologies (IC_ASET)*, Hammamet, Tunisia, 2017.
- [24] H. Wang, Y. Xu and M. Li, "Study on the MFCC similarity-based voice activity detection algorithm," in *2011 2nd International Conference on Artificial Intelligence, Management Science and Electronic Commerce (AIMSEC)*, Deng Feng, China, 2011.
- [25] G.-D. Wu and C.-T. Lin, "Word boundary detection with mel-scale frequency bank in noisy environment," *IEEE Transactions on Speech and Audio Processing*, vol. 8, no. 5, pp. 541-554, 2000.
- [26] B. S. Atal, "Effectiveness of linear prediction characteristics of the speech wave for automatic speaker identification and verification," *The Journal of the Acoustical Society of America*, vol. 55, no. 6, pp. 1304-1312, 1974.
- [27] L. R. Rabiner and R. W. Schafer, *Introduction to Digital Speech Processing*, Now Publishers, 2007.
- [28] E. Zwicke and E. Terhardt, "Analytical expressions for critical-band rate and critical bandwidth as a function of frequency," *The Journal of the Acoustical Society of America*, vol. 68, no. 5, pp. 1523-1525, 1998.
- [29] E. Zwicker, "Subdivision of the Audible Frequency Range into Critical Bands (Frequenzgruppen)," *The Journal of the Acoustical Society of America*, vol. 33, no. 2, pp. 248-248, 1960.
- [30] H. Hermansky, "Perceptual linear predictive (PLP) analysis of speech," *The Journal of the Acoustical Society of America*, vol. 87, no. 4, pp. 1738-1752, 1989.
- [31] A. Kalia, S. Sharma and S. Pandey, "Comparative analysis of speaker recognition system based on voice activity detection technique, MFCC and PLP Features," in *Intelligent Computing Techniques for Smart Energy Systems*, Singapore, 2020.
- [32] "Wavelets in speech segmentation," in *MELECON 2008 - The 14th IEEE Mediterranean Electrotechnical Conference*, Ajaccio, France, 2008.
- [33] T. Ro, T. Marwala and S. Chakraverty, "Precise detection of speech endpoints dynamically: A wavelet convolution based approach," *Communications in Nonlinear Science and Numerical Simulation*, vol. 67, pp. 162-175, 2019.
- [34] C. Yali, L. Dongsheng, J. Shuo and N. Xuefen, "A speech endpoint detection algorithm based on wavelet transforms," in *The 26th Chinese Control and Decision Conference (2014 CCDC)*, Changsha, China, 2014.
- [35] M. Eshaghi and M. .. Mollaei, "Voice activity detection based on using wavelet packet," *Digital Signal Processing*, vol. 20, no. 4, pp. 1102-1115, 2010.

- [36] E. P. Wigner, "On the Quantum Correction for Thermodynamic Equilibrium," in *Part I: Physical Chemistry. Part II: Solid State Physics*, Springer, 1997, pp. 110-120.
- [37] J. Ville, "Theorie et application de la notion de signal analytique," *Câbles et transmissions*, vol. 2, no. 1, pp. 61-74, 1948.
- [۳۸] م.ع. نژادی و ح. حسن پور، "تحلیل و پردازش سیگنال به کمک روش های زمان فرکانس." *مجله علوم رایانشی*, جلد ۲، ۱۳۹۸.
- [39] R. B. Pachori and P. Sircar, "A new technique to reduce cross terms in the Wigner distribution," *Digital Signal Processing*, vol. 17, no. 2, pp. 466-474, 2007.
- [40] F. Auger and P. Flandrin, "Time – Frequency Toolbox For Use with MATLAB," 1995, pp. 1-150.
- [41] S. Group, "SPEECH DATABASE DEVELOPMENT AT MIT: TIMIT AND BEYOND," *Speech Communication*, vol. 2, no. 4, pp. 351-356, 1990.
- [42] R. Subramanian, "Automatic Word Boundary Detection Using Adaptive Scaling", Northeastern University, Boston, MA, USA, 2010.

Abstract

Word boundary detection means the beginning and end of a word in a speech signal, which plays an important role for most applications such as speech recognition, text-to-speech system, and so on. In speech recognition system, most errors are due to incorrect recognition of word boundaries. Correct word boundaries reduces detection errors and improves systems performance. There are many algorithms used for detecting the boundary of words, some of which are not good enough to correct the word boundary due to the noise of the environment or a good pause between words. Which can be reduced by noise methods to be able to identify the borders well. While some methods are combined, they can better detection word boundaries.

In this research, the aim is to detection word boundaries using time-frequency methods such as wavelet and Wigner-Ville. Time-frequency methods examine the signal in both time and frequency domains. Because the speech signal is an unstable signal and its characteristics change with time, the best way to analyze this signal is to extract features from those time-frequency methods. Different methods of extracting MFCC, PLP, LPC, wavelet analysis and Wigner-Ville method were examined to identify word boundaries. Two types of data were used to evaluate the results. In the TIMIT database, time between words is very short and words are pronounced naturally and almost quickly, while in CADLab words are expressed numerically. The Wigner-Will method is able to identify word boundaries with 83.5% for the first data, and 68% for the second data is the best way to identify word boundaries. Entropy with 25% detection for the first data and ZRC with 30% detection for the second data has the lowest rate of correct detection of word boundaries among the methods performed.

Keywords: Word Boundary detection, Feature Extraction, Time-Frequency, MFCC, PLP, LPC



Shahrood University of Technology
Faculty of Electrical and Robatic Engineering

Rock Mechanics MSc Thesis in

**Words boundary detection in speech signal using time-frequency
method**

By:

Fateme Tajik Ijdan

Supervisor:

Dr Hossein Marvi

October 2021