

# پایه گذاری بستری نو و کارآمد در حوزه بازشناسی گفتار فارسی

باقر باباعلی

دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران، تهران، ایران



## چکیده

برخلاف پیشینه سی ساله پژوهش در حوزه بازشناسی گفتار فارسی در ایران و دست یافتن به پیشرفت‌های درخور توجه، نتایج عمده کارهای انجام شده به دلیل عدم وجود بستر یکسان، قابل مقایسه و ارزیابی دقیق نیستند. بستر بیش تر شامل سامانه بازشناسی و دادگان با تعریف مشخص مجموعه‌های آموزش، توسعه و ارزیابی است. سامانه متن‌باز کلدی با وجود نوظهور بودن آن ویژگی‌های منحصر به فردی دارد که در سال‌های اخیر مورد توجه اکثر آزمایشگاه‌های تراز نخست پردازش گفتار دنیا قرار گرفته است و با لحاظ همه جوانب، بهترین انتخاب موجود در راستای پایه‌گذاری این بستر برای تمامی زبان‌ها از جمله زبان فارسی است. در این مقاله پس از بررسی خصوصیات، توانمندی‌ها و اجزای مختلف نرم‌افزار کلدی؛ دادگان فارسی‌دات را به دلیل ثبت رسمی و قابل دسترس بودن آن برای همگان از سراسر دنیا به عنوان بخش دیگر این بستر انتخاب کرده و به ناسی از انتخاب انجام شده بر روی دادگان TIMIT به تعریف مجموعه‌های آموزش، توسعه و ارزیابی می‌پردازیم. در نهایت بیش تر قریب به اتفاق تکنیک‌ها و روش‌های موجود در کلدی بر روی دادگان فارسی‌دات، مطابق تعریف صورت گرفته، مورد آزمایش قرار گرفته‌اند. بهترین میزان خطای حاصل در بازشناسی واج برای مجموعه توسعه ۲۰/۳ درصد و برای مجموعه آزمون ۱۹/۸ بوده است. دسترسی به کدهای نوشته در جهت فراهم‌سازی این بستر، در نرم‌افزار کلدی موجود است که با توجه به متن‌باز بودن آن، دسترسی به آنها به منظور بازسازی نتایج آمده در این مقاله در صورت در اختیار داشتن دادگان فارسی‌دات به راحتی قابل انجام است.

واژگان کلیدی: بازشناسی گفتار پیوسته فارسی، دادگان فارسی‌دات، نرم‌افزار متن‌باز کلدی.

## A state-of-the-art and efficient framework for Persian Speech Recognition

Bagher BabaAli

School of Mathematics, Statistics and Computer Sciences, College of Science, University of Tehran, Tehran, Iran

### Abstract

Although researches in the field of Persian speech recognition claim a thirty-year-old history in Iran which has achieved considerable progresses, due to the lack of well-defined experimental framework, outcomes from many of these researches are not comparable to each other and their accurate assessment won't be possible. The experimental framework includes ASR toolkit and speech database which consists of training, development and test datasets. In recent years, as a state-of-the-art open-source ASR toolkit; Kaldi has been very well-received and welcomed in the community of the world-ranked speech researchers

\* نویسنده عهده‌دار مکاتبات

سال ۱۳۹۵ شماره ۳ پیاپی ۲۹

around the world. considering all aspects, Kaldi is the best option among all of the other ASR toolkits to establish a framework to do research in all languages, including Persian.

In this paper, we chose Fardat as the speech database which is the counterpart of TIMIT for Persian language because not only it has got a standard form but it's also accessible for all researchers around the world. Similar to the recipe on TIMIT database, we defined these three sets on the Farsdat: Training, Development and Test sets. After a survey on Kaldi's components and features, we applied most of state-of-the-art ASR techniques in the Kaldi on the Farsdat based on three sets definition. The best phone error rate on development and test set have been 20.3% and 19.8%. All of the codes and the recipe that was written by author have been submitted to Kaldi repository and they are accessible for free, so all the reported results will be easily replicable if you have access to Farsdat database.

**Keyword:** Persian Continuous Speech Recognition, FarsDat Database, Kaldi Toolkit

جانز هاپکینز؛ هسته اصلی یک سامانه بازشناسی گفتاری پیوسته‌ای شکل گرفت که کلدی<sup>۱</sup> [8] نام گرفت. برخلاف این که از زمان پایه‌عرصه‌گذاشتن کلدی زمان زیادی نمی‌گذرد، ولی در بیت پژوهش‌گران علوم گفتار دنیا به دلیل ویژگی‌های منحصربه‌فرد آن به شدت مورد توجه قرار گرفته است. تمامی تکنیک‌های کارآمد و به‌روز بازشناسی گفتار یا در این سامانه پیاده‌سازی شده و یا در حال پیاده‌سازی است. تا این حد که شبکه‌های عصبی ژرف<sup>۲</sup> [9]، [10] که داغ‌ترین مبحث پژوهشی در سال‌های اخیر در این عرصه است، برای مدل‌سازی آکوستیکی [11]، [12] و همچنین امتیازدهی مجدد<sup>۳</sup> امتیاز مدل زبانی [13] پیاده‌سازی شده است و پیوسته در حال توسعه و بهبود است. از حدود دو دهه قبل در عرصه بازشناسی گفتار، سامانه‌های متن‌باز متعددی نظیر HTK [14]، SONIC [15]، [16]، SPHINX [17]، [18]، RASR [19]، [20]، JULIUS [21]، SPRAAK [22]، KALDI و در همین‌واخر [23] عرضه شده است. با لحاظ همه جوانب و معیارها؛ تنها رقیب سامانه کلدی، سیستم RASR محسوب می‌شود که توسط آزمایشگاه RWTH دانشگاه آخن آلمان توسعه داده شده است. کلدی در مقایسه با RASR از معماری و ساختار نرم‌افزاری به مراتب بهتری برخوردار است که کار فهم کد، اعمال تغییرات و توسعه را برای پژوهش‌گران تسهیل می‌کند. و علاوه بر موارد ذکر شده، محدودیت قانونی به‌کارگیری در کاربردهای تجاری ندارد و عزم و ارائه توسعه‌دهندگان آن بر پیاده‌سازی آخرین تکنیک‌های کارآمد بازشناسی گفتار و به‌روز نگه‌داشتن آن راسخ‌تر بوده است.

با عنایت به برتری سامانه متن‌باز کلدی از جمیع جوانب و نیاز به وجود یک سامانه بازشناسی گفتار متن‌باز و

<sup>1</sup> Kaldi

<sup>2</sup> Deep Neural Network

<sup>3</sup> Rescoring

## ۱- مقدمه

به‌منظور تسهیل ارتباط انسان با ماشین، پژوهش در حوزه بازشناسی گفتار از حدود نیم‌قرن، همواره مورد توجه پژوهش‌گران علوم برق و رایانه بوده است. با توجه وابستگی این مقوله به زبان؛ لذا با تاسی از این پژوهش‌ها در دو دهه اخیر در حوزه بازشناسی گفتار فارسی پژوهش‌های متعددی توسط پژوهش‌گران داخلی این عرصه صورت گرفته است [1]-[7].

با وجود دستیابی این پژوهش‌گران به نتایج قابل قبول و ارائه یک‌سری محصولات تجاری بر مبنای آن؛ ولی از دیدگاه پژوهشی همچنان مشکلات متعددی پابرجاست که در ادامه به آن می‌پردازیم. پژوهش‌های انجام‌شده بیش‌تر بر مبنای سامانه‌های بازشناسی گفتار و دادگان‌های متنوعی گزارش شده است که ارزیابی و مقایسه نتایج دست‌یافته را ناممکن می‌کند. برخی از این سامانه‌ها در آزمایشگاه خاصی توسعه داده شده و دسترسی به آن برای دیگران مقدور نیست. برخی از این پژوهش‌ها بر مبنای دادگان‌هایی بوده که چندان در دسترس همگان نبوده و یا استاندارد نیست. به این ترتیب مقایسه نتایج و ارزیابی تکنیک‌های ارائه شده ناممکن است.

بنابراین شکی نیست که در مقوله بازشناسی گفتار، فراهم‌آوردن شرایطی که قضاوت در مورد روش‌ها و تکنیک‌های ارائه شده توسط پژوهش‌گران را تسهیل کند یک نیاز اساسی است. برای تحقق این امر مهم می‌بایست یک سامانه بازشناسی گفتار پایه که به‌روز هم است، در دسترس همگان باشد. همچنین نتایج بر مبنای دادگان‌های استاندارد با شرایط مشابه ارائه شود که قضاوت را به مراتب روشن و آسان‌تر می‌کند. در راستای تحقق این امر، در کارگاه تابستانی سال ۲۰۰۹ در مرکز پژوهش زبان و گفتار دانشگاه

در ادامه، در بخش بعد به بررسی سامانه کلدی شامل روش‌ها، تکنیک‌های و قابلیت‌های موجود در آن می‌پردازیم. بخش سوم به بررسی دادگان فارس‌دات و شرایط آزمایش تعریف‌شده بر روی آن اختصاص داده شده است. نتایج و بررسی و تحلیل آنها در بخش چهارم آمده و در نهایت در بخش پنجم جمع‌بندی ارائه شده است.

## ۲- سامانه بازشناسی گفتار کلدی

کلدی یک سامانه متن باز برای بازشناسی گفتار است که در محیط در لینوکس به زبان ++C توسعه داده شده است. هرچند که قابلیت به‌کارگیری در محیط ویندوز را هم دارد. خلق کلدی با هدف ارائه یک کد مدرن، انعطاف‌پذیر و ساده جهت فهم، تغییر و توسعه صورت گرفت. از ویژگی‌های مهم کلدی می‌توان این موارد را برشمرد: برای پیاده‌سازی بخش کدگشایی<sup>۴</sup> به روش ایستای مبنی بر FST<sup>۵</sup> [27] که به‌روزترین و کارآمدترین روش در این حوزه است، عمل شده است و برای پیاده‌سازی FST از کتابخانه متن باز OpenFST [28] بهره گرفته شده است. در روش‌های ایستا کدگشایی برخلاف روش‌های پویا؛ کل فضای جستجو به صورت یک گراف بزرگ، برون‌خط ساخته شده و در هنگام بازشناسی در آن فقط جستجو صورت می‌گیرد. اندازهٔ گراف حاصل به اندازهٔ مدل زبانی و مدل آکوستیکی و تعداد مدخل واژگان بستگی دارد. در کلدی برای انجام محاسبات ماتریسی؛ از یک طریق یک واسط نرم‌افزاری، توابع کتابخانه‌های استاتارد LAPACK و BLAS فراخوانی می‌شوند. الگوریتم‌ها در کلدی به عمومی‌ترین شکل ممکن طراحی شده‌اند، به‌نحوی که توسعه به‌راحتی قابل انجام باشد. برای مثال کدگشا در کلدی از طریق یک واسط نرم‌افزاری امتیاز هر قاب گفتار را به‌ازای هر حالت HMM دریافت می‌کند. این امتیاز می‌تواند هم بر مبنای SGMM<sup>۱</sup> [29] یا GMM محاسبه شده باشد و یا از خروجی شبکه عصبی باشد. کلدی تحت لیسانس Apache v2,0 عرضه شده که یکی از کم‌محدودیت‌ترین‌هاست. در کلدی دستورالعمل کامل برای ساخت یک سامانه بازشناسی گفتار بر مبنای اکثر قریب به اتفاق دادگان‌های گفتاری استاندارد در قالب یک فایل bash ارائه شده است که کمک بزرگی به پژوهش‌گران محسوب می‌شود. برای بیش‌تر کدهای کلدی توبلجی جهت تست صحت عملکرد آنها نوشته شده است. هر چند در حال حاضر

به‌روز برای انجام پژوهش‌ها در حوزه بازشناسی گفتار فارسی، بر آن شدیم که در سامانه کلدی، بستری برای این منظور بنا کنیم. از بین دادگان‌های گفتاری موجود برای زبان فارسی دادگان فارس‌دات [24]، استاندارد، ثبت شده و در سطح بین‌المللی عرضه شده است که تهیه آن را برای پژوهش‌گران در سراسر دنیا مقدور می‌کند. این دادگان توسط پژوهشگاه توسعه فناوری‌های پیشرفتهٔ خواجه نصیرالدین طوسی برای اهداف متعددی در حوزه زبان‌شناسی و پردازش گفتار زبان فارسی طراحی و گردآوری شده است و شباهت زیادی به دادگان انگلیسی TIMIT [25] دارد. به‌منظور تعادل فنوتیکی، جملات آن طوری انتخاب شده که پوشش خوبی بر روی واج‌ها و دوواچی‌های زبان فارسی داشته باشد. بنابراین متن برخی جملات ساختگی، یعنی چندان متداول نیست. با توجه به ساختگی بودن جملات و کوچک‌بودن حجم؛ این دادگان مشابه دادگان TIMIT برای بازشناسی واج مناسب است نه کلمه. بنابراین دادگان فارس‌دات جهت بنای بستر بازشناسی گفتار فارسی در سامانه کلدی انتخاب شد.

به‌طورمعمول در ارزیابی‌های متداول در حوزهٔ بازشناسی گفتار سه مجموعه آموزش<sup>۱</sup>، توسعه<sup>۲</sup> و ارزیابی<sup>۳</sup> یا آزمون در نظر گرفته می‌شود. بر مبنای مجموعه آموزش؛ مدل‌های آکوستیکی واحدهای آوایی آموزش داده می‌شود. مجموعه توسعه برای تنظیم و بهینه‌سازی پارامترها از جمله وزن مدل زبانی استفاده می‌شود و در آخر براساس مدل‌های آموزش داده شده و پارامترهای تنظیم‌شده، مجموعه آزمون مورد آزمایش قرار می‌گیرد. با توجه به چندمنظوره بودن دادگان فارس‌دات، این سه مجموعه برای این دادگان تعریف نشده است. در انتخاب این سه مجموعه برای دادگان TIMIT سعی شده که هیچ‌گونه گوینده مشترکی بین آنها نباشد و همچنین بین متن جملات ادا شده توسط گویندگان مجموعه آموزش از یک طرف و مجموعه توسعه و آزمون از طرف دیگر جملهٔ یکسانی وجود ندارد [26]. در این مقاله نیز با لحاظ این دو محدودیت ذکرشده؛ بر روی دادگان فارس‌دات سه مجموعه ذکر شده تعریف شده و در آزمایش‌های انجام‌شده، بر مبنای آن عمل شده است. در این مقاله هدف تعریف یک بستر استاندارد و به‌روز در سامانه کلدی برای پژوهش‌گران بازشناسی گفتار فارسی بر روی دادگان فارس‌دات و ارائه نتایج حاصله به‌منظور انجام مقایسات بعدی است.

<sup>4</sup> Decoding

<sup>5</sup> Finite State Transducer

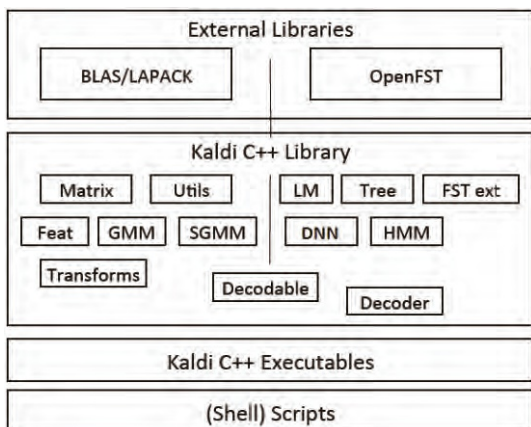
<sup>6</sup> Sub-space Gaussian Mixture model

<sup>1</sup> Training Set

<sup>2</sup> Development Set

<sup>3</sup> Evaluation Set

آن حالت‌های تمام HMMهاست به دست می‌آید که الحاق آن با بردارهای ویژگی متداول، بهبود درخور توجهی را در دقت بازشناسی به همراه دارد. در نسخه برخط کلدی که در همین اواخر پیاده‌سازی شده‌است، به دلیل آن امکان به‌کارگیری مؤثر<sup>۵</sup> fMLLR [32] و CMVN<sup>۶</sup> [33] وجود ندارد، از بردار ویژگی ivector [34] بهره گرفته شده که به‌طور کامل افت دقت ناشی از عدم به‌کارگیری دو تکنیک مؤثر ذکر شده را جبران می‌کند، به‌نحوی که هیچ‌گونه تفاوتی در دقت بین نسخه برخط و برون‌خط کلدی وجود ندارد [12].



(شکل - ۱): نمایی ساده‌شده از اجزای نرم‌افزاری سامانه کلدی  
(Figure-1): A simplified diagram on the software modules from the Kaldi toolkit

در روش استخراج ویژگی ivector که در سالیان اخیر بیش‌تر برای بازشناسی گوینده ارائه شده است، از گفتار با طول متغیر یک بردار ویژگی با طول ثابت استخراج می‌کند که گویای خصوصیات گوینده و مستقل از متن گفتار است. با الحاق این بردار ویژگی به بردار ویژگی MFCC یا PLP و به‌کارگیری آن در بازشناسی گفتار، اثر عدم اعمال fMLLR و CVN در حالت به‌کارگیری برخط جبران می‌شود. در بخش استخراج ویژگی کلدی، امکان به‌کارگیری VTLN [35]، [36]، به‌منظور هنجارسازی اثر طول کانال صوتی گوینده و CMVN به‌منظور حذف اثر کانال و اعمال تبدیلات خطی LDA<sup>۷</sup>، HLDA<sup>۸</sup> [37] و MLLT<sup>۸</sup> [38] وجود دارد. از آنجا که به‌کارگیری فرکانس گام در بهبود دقت بازشناسی زبان‌های آهنگین نظیر چینی و ویتنامی دارد [39]، در کلدی یک روش استخراج فرکانس گام، اختصاصی‌شده به‌منظور

کلدی از روش‌های کارآمد و به‌روزی برای مدل‌سازی زبانی و کدگشایی بهره برده است و قابلیت کدگشایی برخط را هم دارد، ولی در ابتدا بیشتر با هدف بستری برای انجام پژوهش‌ها در حوزه مدل‌سازی آکوستیکی عرضه شد. در ادامه در این بخش به تشریح ساختار کلدی، توصیف مؤلفه‌های اصلی بازشناسی گفتار در کلدی نظیر بخش استخراج ویژگی، مدل‌سازی آکوستیکی، درخت تصمیم فونتیکی، مدل‌سازی زبانی و در آخر کدگشا خواهیم پرداخت.

## ۱-۲- ساختار نرم‌افزاری

در شکل (۱) ساختار نرم‌افزاری کلدی آمده است. همان‌طوری که در این شکل مشاهده می‌شود در آن از دو کتابخانه خارجی متن‌باز استفاده است: یکی OpenFST و دیگری BLAS<sup>۱</sup>/LAPACK<sup>۲</sup> که به‌منظور انجام محاسبات عددی به کار گرفته شده است. دیگر کتابخانه‌های کلدی که در شکل نمایش داده به‌صورت داخلی هستند. این کتابخانه‌ها به‌صورت طبقه در زبان C++ نوشته شده‌اند و فراخوانی توابع آنها از طریق برنامه‌های اجرای که به‌صورت خط فرمان اجرا می‌شود، صورت می‌گیرد. برای ایجاد یک سامانه بازشناسی، این توابع اجرایی در یک زبان اسکریپتی مانند bash فراخوانی می‌شوند. هر یک از این توابع اجرایی یک وظیفه مشخص برای آن تعریف شده است که از طریق آرگومان‌های خط فرمان به‌نحو دقیق‌تری مشخص می‌شود. همه این توابع اجرایی می‌توانند عمل خواندن از و یا نوشتن به خروجی را از طریق لوله<sup>۳</sup> انجام دهند که این امکان به‌رشته درآوردن تعدادی از آنها را فراهم می‌کند.

## ۲-۲- استخراج ویژگی

سه روش استخراج ویژگی MFCC، PLP [30] و BNF<sup>۴</sup> [31] موجود است و پارامترهای هر یک از این روش‌ها دارای مقدار پیش فرضی می‌باشد، هرچند که امکان تغییر و تنظیم آن توسط کاربر در قالب فایل Config وجود دارد. روش استخراج ویژگی BNF در سال‌های اخیر ارائه و مورد توجه قرار گرفته است. در این روش، بردار ویژگی از لایه<sup>۴</sup> میانی یک شبکه عصبی ژرف که ورودی آن بردار ویژگی فعلی به همراه تعداد محدودی بردارهای قبل و بعد (سه یا پنج) و خروجی

<sup>۵</sup> Feature-space Maximum Likelihood Linear Regression

<sup>۶</sup> Cepstral Mean and Variance Normalization

<sup>۷</sup> Heteroscedastic Linear Discriminant Analysis

<sup>۸</sup> Maximum Likelihood Linear Transform

<sup>۱</sup> Basic Linear Algebra Subroutines

<sup>۲</sup> Linear Algebra PACKage

<sup>۳</sup> Pipe

<sup>۴</sup> Bottleneck Feature

۲-۳-۲- مدل‌های مخلوط گوسی زیرفضا (SGMM)  
این روش مدل‌سازی آکوستیکی با الهام از بحث تحلیل مؤلفه‌های مشترک<sup>۵</sup> که یک رویکرد جدید در مبحث بازشناسی گوینده محسوب می‌شود، در کارگاه تابستانی سال ۲۰۰۹ مرکز زبان و گفتار دانشگاه جانزهاپکینز ارائه شد. به‌طورکلی در این روش یک مدل مخلوط جهانی (UBM) بر مبنای کل گفتار آموزش داده می‌شود و آنگاه با تطبیق آن به مدل‌های مربوط به هر حالت می‌رسیم. در مواقعی که داده آموزشی در محدوده<sup>۶</sup> کمی در اختیار است، به‌عنوان مثال بازشناسی گفتار برای زبان‌ها با منابع محدود<sup>۷</sup>، این روش امکان به‌رمندی از داده خارج از محدوده<sup>۸</sup> و یا خارج از زبان<sup>۹</sup> را جهت جبران کمبود داده در این موارد فراهم می‌کند [29].

### ۲-۳-۳- شبکه‌های عصبی ژرف (DNN)

به‌کارگیری شبکه‌ی عصبی به‌صورت ترکیب با HMM (به جای GMM) برای مدل‌سازی آکوستیکی به‌دلیل قدرت تمایز بالای آن در مقایسه با GMM در دهه نود میلادی مورد توجه پژوهش‌گران بازشناسی قرار گرفت [41]، [42]. در روش ترکیبی، یک شبکه‌ی عصبی جهت تخمین احتمال پسین<sup>۱۰</sup> هر بردار ویژگی به‌ازای حالات HMMها آموزش داده می‌شود که درنهایت با تقسیم آنها بر احتمال پیشین<sup>۱۱</sup> هر حالت؛ به میزان شباهت<sup>۱۲</sup> حالات می‌رسیم. در تلاش‌های نخستین که در دهه نود صورت گرفت، به‌دلیل محدودبودن توان محاسباتی، تعداد لایه‌های مخفی شبکه‌ی عصبی یک و یا حداکثر دو لایه بود و هیچ‌گونه تلاشی در جهت مقداردهی نخستین وزن‌های شبکه‌ی عصبی به‌صورت هوشمندانه صورت نمی‌گرفت. موفقیت رویکرد ترکیبی [43]، [45]، [10]، [12]، در سال‌های اخیر بیش‌تر مدیون بکارگیری شبکه‌های عصبی ژرف (تعداد لایه‌های مخفی بیشتر) به لطف افزایش توان محاسباتی رایانه‌ها و به‌کارگیری ماشین بولتزمن محدودیت‌دار (RBM<sup>۱۳</sup>) [46] برای پیش‌آموزش شبکه‌ی عصبی بوده است. در کلدی دو کد برای پیاده‌سازی سامانه ترکیبی به‌صورت مجزا توسط Daniel Povey و Karel Vesely توسعه داده شده است. عمده تمایز این دو سامانه در

بازشناسی گفتار ارائه شده است. در روش ارائه‌شده [40] به‌ازای هر قاب گفتار دو مؤلفه<sup>۱</sup> POV و فرکانس گام که در بخش‌های بی‌صدا مقدار آن درون‌یابی شده؛ استخراج و به سایر ویژگی‌های استخراج‌شده از آن قاب گفتار الحاق می‌شود.

روند مرسوم در کلدی برای استخراج ویژگی به این صورت است: در آموزش نخستین سامانه بر مبنای تک‌واج<sup>۲</sup>، بردار ویژگی MFCC یا PLP به همراه مشتقات نخست و دوم آن به‌صورت یک بردار ۳۹ مؤلفه‌ای استفاده می‌شود. در آموزش و ارزیابی سامانه بر مبنای سه‌واجی<sup>۳</sup> بر روی هر بردار ویژگی ۱۳ مؤلفه‌ای به همراه سه بردار ویژگی قبل و بعد تبدیل LDA و سپس تبدیل MLLT اعمال شده و درنهایت به یک بردار ویژگی چهار مؤلفه‌ای تبدیل می‌شود. با توجه به خطی‌بودن این تبدیلات، ماتریس آنها در مرحله آموزش و بر مبنای داده‌های آموزشی محاسبه می‌شود.

### ۲-۳-۲- مدل‌سازی آکوستیکی

در کلدی سه روش مدل‌سازی آکوستیکی GMMs، SGMMs و DNN موجود است و با توجه به تدابیر اندیشیده شده در معماری کلدی، به‌راحتی می‌توان سایر روش‌های مدل‌سازی را هم پیاده‌سازی و به‌کار گرفت.

### ۲-۳-۱- مدل‌های مخلوط گوسی (GMM)

مدل‌های مخلوط گوسی با ماتریس کواریانس کامل و قطری در کلدی موجود است و پیاده‌سازی آنها به‌نحوی صورت گرفته که محاسبه میزان شباهت به‌ازای هر GMM تنها با ضرب نقطه‌ای قابل انجام است که افزایش سرعت زیادی را به همراه دارد. در کلدی امکان توصیف مجزای توپولوژی مدل مخفی مارکوف مربوط به هر واج وجود دارد. توپولوژی هر مدل مخفی مارکوف شامل تعداد حالت‌ها و نحوه پرش بین آنها است. در کلدی مقدار پیش‌فرض تعداد حالت‌ها برای واج‌ها سه حالت است و برای سکوت و نوفه و سایر اصوات غیرکلامی پنج حالت در نظر گرفته شده است. بطور واضح مدل HMM واج‌ها چپ به راست و مدل HMM سکوت و نوفه و اصوات غیرکلامی ارگودیک<sup>۴</sup> است؛ یعنی از هر حالت به تمام حالات دیگر امکان پرش وجود دارد.

<sup>5</sup> Joint Factor Analysis

<sup>6</sup> In-domain

<sup>7</sup> Low Resource Languages

<sup>8</sup> Out-of-domain

<sup>9</sup> Out-of-language

<sup>10</sup> Posterior Probability

<sup>11</sup> Prior Probability

<sup>12</sup> Likelihood

<sup>13</sup> Restricted Boltzmann Machine

<sup>1</sup> Probability Of Voicing

<sup>2</sup> Monophone

<sup>3</sup> Triphone

<sup>4</sup> Ergodic

گسترش یافته می‌شود. در این موارد ریشه‌های درخت‌های تصمیم بین صورت‌های مختلف یک واج مشترک است.

## ۲-۴- مدل‌سازی زبانی

از آنجا که کدگشای کلیدی بر مبنای روش استاتیک FST عمل می‌کند، بنابراین قادر به استفاده از هر گونه مدل زبانی است که به یک FST قابل تبدیل باشد. ابزاری جهت تبدیل مدل زبانی از فرمت ARPA به شکل FST وجود دارد. در مواقعی که نیاز به هرس مدل زبانی<sup>۷</sup> باشد از نرم‌افزار IRSTLM [50] استفاده می‌شود. برای ساخت مدل زبانی از روی متن خام نیز همین نرم‌افزار IRSTLM قابل استفاده است؛ هرچند که می‌توان از نرم‌افزار قدرتمند SRILM [51] هم استفاده کرد.

به‌تازگی استفاده از شبکه‌های عصبی بازگشتی<sup>۸</sup> برای مدل‌سازی زبانی مورد توجه پژوهش‌گران قرار گرفته است [13]. در کلیدی این ابزار مدل‌سازی زبانی که به‌طور مجزا توسعه داده شده، موجود است و امکان به‌کارگیری آن فراهم شده است. یک روش متداول برای به‌کارگیری آن، امتیازدهی مجدد امتیاز مدل زبانی خروجی‌های بازشناسی (N-best) است که به‌طور معمول بهبود دقت بازشناسی را به همراه دارد.

## ۲-۵- کدگشا

در کلیدی مرحله بازشناسی گفتار از WFST استفاده می‌شود. فرایند ساخت گراف کدگشایی بر مبنای روش آمده در [27] است هر چند که تغییراتی به‌منظور بهبود انجام شده است. به‌عنوان مثال این تغییر صورت گرفته که FST تصادفی<sup>۹</sup> باشد؛ به این معنی که مجموع اوزان خروجی‌ها در هر حالت برابر یک باشد. مدل‌های HMM واج‌ها (H)، وابستگی به متن (C)، واژگان (L) و مدل زبانی (G) به‌صورت FST مجزا با هم ترکیب شده و گراف بزرگی تحت عنوان HCLG را می‌سازند که فضای جستجو را در هنگام کدگشایی شکل می‌دهد.

در کلیدی، بسته به برخط یا برون‌خط بودن فرایند بازشناسی، نوع خروجی (شبکه<sup>۱۰</sup> یا n-تا بهترین خروجی)، میزان بهینه‌بودن و سرعت اجرا، کدگشاهای متعددی وجود دارد. همگی کدگشاها این ویژگی را دارند که به‌راحتی

آن است که Karel مقدردهی نخستین شبکه را با روش پیش‌آموزش RBM و Daniel به‌صورت تصادفی انجام می‌دهد. این دو سامانه به‌لحاظ دقت بازشناسی، بسیار مشابه هم بوده و در برخی مواقع سامانه Karel به میزان کمی بهتر عمل می‌کند. با توجه به این که به‌کارگیری شبکه عصبی ژرف در بازشناسی گفتار، داغ‌ترین مبحث پژوهشی این حوزه در سالیان اخیر بوده است، در سامانه کلیدی پیوسته روش‌های جدیدی در این حوزه توسط پژوهش‌گران سراسر دنیا توسعه داده شده و آزمایش می‌شود که پرداختن به آن خارج از حوصله این مقاله است.

## ۲-۳-۴- تطبیق گوینده

تطبیق گوینده<sup>۱</sup> در کلیدی به دو صورت، یکی در فضای مدل (MLLR) و دیگری در فضای ویژگی (fMLLR) قابل انجام است. برای هر دو روش می‌بایست تبدیلات لازم با استفاده از درخت رگرسیون تخمین زده شود [47]. همچنین امکان هنجارسازی گوینده با استفاده از VTLN و یا به‌صورت کلی‌تر هنجارسازی جنسیت<sup>۲</sup> با استفاده از تبدیل‌نمایی<sup>۳</sup> [48] وجود دارد. برای آموزش تطبیق به گوینده (SAT<sup>۴</sup>) مدل‌های آکوستیکی، امکان استفاده از هر دو روش VTLN و fMLLR وجود دارد که استفاده از fMLLR متداول‌تر است.

## ۲-۳-۵- درخت‌های تصمیم فونتیکی

در کلیدی کد ساخت درخت تصمیم فونتیکی به‌نحوی است که امکان مدل‌سازی بافت<sup>۵</sup> با هر طولی را مقدور می‌سازد و همچنین به‌اندازه کافی عمومی است که امکان به‌کارگیری طیف وسیعی از روش‌ها را دارد. روش کلاسیک به این صورت است که به ازای هر حالت HMM مربوط به هر تک‌واج یک درخت تصمیم وجود دارد، که بر مبنای سؤالاتی در مورد واج‌های سمت چپ و راست تصمیم‌گیری می‌کند. در کلیدی، ریشه‌های درخت‌های تصمیم بین واج‌ها و حالت‌های یک واج می‌تواند به اشتراک گذاشته شود. سؤالات فونتیکی بر مبنای دانش زبان‌شناسی می‌تواند تهیه شود؛ ولی در کلیدی سؤالات بر مبنای خوشه‌بندی درختی<sup>۶</sup> واج‌ها به‌صورت خودکار تولید می‌شوند [49]. سؤالات در مورد چیزهایی نظیر استرس فونتیکی - در صورتی که در واژگان آمده باشد- و اطلاعات شروع و پایان کلمات منجر به یک مجموعه واج

<sup>1</sup> Speaker Adaptation

<sup>2</sup> Gender Normalization

<sup>3</sup> Exponential Transform

<sup>4</sup> Speaker Adaptive Training

<sup>5</sup> Context

<sup>6</sup> Tree-Clustering

<sup>7</sup> Language Model Pruning

<sup>8</sup> Recurrent Neural Network

<sup>9</sup> Stochastic

<sup>10</sup> Lattice

بررسی دقیق این دادگان و الهام از آنچه برای دادگان TIMIT در [26] صورت گرفته است، سعی شد این سه مجموعه را انتخاب کنیم. با لحاظ این محدودیت که بین سه مجموعه، گوینده مشترکی نباشد و متن جملات مجموعه آموزش با مجموعه توسعه و آزمون اشتراکی نداشته باشد، این انتخاب صورت گرفت. همچنین سعی شده تا حد امکان توزیع نوع گویش، نسبت مرد به زن، میزان تحصیلات و سن بین این سه مجموعه مطابق توزیع کل دادگان باشد که چندان کار راحتی نبود و نیاز به بررسی و استخراج آمارهای متعدد و لحاظ آنها داشت. مشخصات سه مجموعه انتخاب شده، مطابق جدول (۱) است. همان طوری که در قبل ذکر شد، تعداد کل جملات بیان شده در فارسات ۶۰۸۰ است؛ در حالی که مجموع جملات سه مجموعه تعریف شده (۲۸۷+۴۷۵+۳۹۹۴) جمله است. شماره گویندگان هریک از این سه مجموعه نیز در جدول (۲) آمده است. برخی از جملات بیان شده توسط بعضی از گویندگان کنار گذاشته شده است؛ به این دلیل که این محدودیت را رعایت کرده باشیم: هیچ گونه جمله با متن مشابه بین مجموعه آموزش از یک طرف و مجموعه توسعه و آزمون وجود نداشته باشد. در تعریف این سه مجموعه بر روی دادگان TIMIT هم به خاطر رعایت این محدودیت، تعدادی از جملات بیان شده (در همین مقیاس) کنار گذاشته شده است.

(جدول-۱): مشخصات سه مجموعه انتخاب شده بر روی دادگان

فارسات

(Table-1): Specification of three datasets defined on the FarsDat database

تعداد جملات	تعداد گویندگان (نفر)	مجموعه
3994	224	آموزش
475	50	توسعه
287	30	آزمون

#### ۴- نتایج

سامانه کلدی در محیط لینوکس توسعه داده شده است و کدنویسی آن در سه سطح صورت گرفته است. به طور معمول توابع پایه‌ای به زبان برنامه نویسی C++ نوشته می‌شود و این توابع در محیط شل اسکریپت برای انجام کاری خاص نظیر استخراج ویژگی، آموزش، اعمال تبدیلات به منظور مقاوم سازی و کدگشایی فراخوانی می‌شوند. برای دادگان‌های استاندارد متعددی، در قالب یک فایبل شل، کارهای

می‌توان نوع مدل آکوستیکی آنها را تغییر داد. کدگشاهای موجود در کلدی در سطح کد C++ همگی تک‌گذره هستند و چندگذره بودن در سطح اسکریپت پیاده‌سازی شده است.

### ۳- دادگان فارسات

فارسات اولین دادگان استاندارد گفتاری زبان فارسی است که با هدف مطالعه مبانی و مدل‌سازی آکوستیکی زبان فارسی به منظور به‌کارگیری در سامانه‌های بازشناسی گفتار فارسی توسط پژوهشگاه توسعه فناوری‌های پیشرفته‌ی خواجه نصیرالدین طوسی تولید شده است [24]. این دادگان شامل ۳۰۴ گوینده فارسی‌زبان با یکی از ده لهجه رایج تهرانی، ترکی، اصفهانی، شمالی، جنوبی، لری، کردی، یزدی، خراسانی و بلوچی با سن، جنسیت و میزان تحصیلات متنوع است. نسبت گویندگان مرد به زن دو به یک است. هر گوینده بیست جمله را در دو جلسه مجزا در یک اتاق آکوستیک به سبک رسمی بیان کرده است و بنابراین تعداد کل جملات ۶۰۸۰ است. جملات به صورت ۱۶ بیتی با فرکانس نمونه‌برداری ۲۲۰۵۰ هرتز و با نسبت سیگنال به نویز<sup>۱</sup> متوسط ۳۲ دسی‌بل ضبط شده است. تعداد کل جملات فارسات ۴۰۵ جمله است که دو جمله آن شامل کل واج‌های زبان فارسی (به جز /f/) است و این دو جمله توسط همه گویندگان بیان شده است. جملات به‌گونه‌ای انتخاب شده‌اند که شامل تمامی دوواچی‌های رایج زبان فارسی باشد. فارسات استاندارد مؤسسه اروپایی داده‌های زبانی<sup>۲</sup> را کسب کرده و در کاتالوگ این مؤسسه ثبت شده است.

در ارزیابی‌های متداول در حوزه بازشناسی گفتار در دنیا، دادگان گفتاری به سه زیرمجموعه آموزش، توسعه و آزمون یا ارزیابی تقسیم می‌شود. بر روی مجموعه آموزش مدل آکوستیکی آموزش داده می‌شود. بر روی مجموعه توسعه، بهینه‌سازی و تنظیم پارامترها انجام می‌شود و در نهایت دقت بر روی مجموعه آموزش با فرض مشخص و ثابت شدن پارامترها بر مبنای مجموعه توسعه گزارش می‌شود. با عنایت به این که فارسات برای اهداف متنوعی گردآوری شده، بر روی آن این سه مجموعه به‌طور رسمی تعریف نشده است؛ و در نتایج گزارش شده، مجموعه توسعه در نظر گرفته نشده است و یا اشتراکی بین این سه مجموعه از نظر گوینده و یا متن جمله وجود دارد. در این پژوهش با

<sup>1</sup> Signal to Noise Ratio

<sup>2</sup> European Linguistic Resource Association (ELRA)

نوشته شده و در اختیار همگان است. با دراختیارداشتن دادگان فارس‌دات و نصب سامانه کلدی، کلیه نتایج آمده در ادامه به راحتی و بدون هیچگونه زحمت برنامه‌نویسی قابل تولید مجدد است. به این ترتیب پژوهش‌گران حوزه بازشناسی گفتار زبان فارسی به راحتی می‌توانند به پیاده‌سازی ایده‌های خود و مقایسه آن در یک بستر بروز بپردازند.

آماده‌سازی و آمایش دادگان، واژگان و مدل زبانی، و فراخوانی فایل‌های شل به منظور استخراج ویژگی، آموزش و کدگشایی صورت گرفته است. برای دادگان فارس‌دات هم این کار توسط نویسنده این مقاله انجام گرفته است. کدهای لازم برای آماده‌سازی سه مجموعه آموزش، توسعه و آزمون تعریف شده بر روی دادگان فارس و انجام آزمایش‌های زیر، در سامانه کلدی توسط نگارنده مقاله،

(جدول-۲): شماره گویندگان سه مجموعه انتخاب شده بر روی دادگان فارس‌دات  
(Table-2): Speaker numbers of three datasets defined on the FarsDat database

مجموعه	شماره گویندگان
آموزش	001 002 003 004 005 006 008 011 012 013 014 015 016 017 019 020 021 022 023 024 025 026 027 028 029 030 031 033 034 035 038 039 040 041 042 043 045 047 048 050 051 052 053 054 055 056 057 058 059 060 061 062 064 065 066 067 068 069 070 071 072 073 074 075 076 077 078 079 081 082 083 084 085 086 087 088 089 091 092 093 094 095 096 099 100 101 102 104 105 106 107 108 109 110 111 112 114 116 117 119 120 121 122 123 124 126 127 128 129 130 131 132 133 134 136 137 138 140 141 142 145 146 147 148 151 153 154 156 157 158 160 161 163 164 165 166 170 171 173 174 175 176 177 178 179 180 182 183 184 185 186 187 188 189 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 207 208 210 211 215 216 217 218 220 221 223 224 227 228 229 230 230 231 232 234 235 237 238 240 241 243 244 245 246 247 248 249 251 252 253 256 257 258 260 261 262 264 265 266 267 268 269 279 293 297 298 300 301 302 303 304
توسعه	009 018 036 046 049 063 080 098 103 113 115 118 135 144 149 150 159 162 168 169 190 206 209 214 222 226 239 254 255 259 263 270 271 272 276 277 278 281 282 284 285 286 287 289 290 291 292 294 296 299
آزمون	007 010 032 037 044 090 097 125 139 143 152 155 167 172 181 212 213 219 225 233 236 242 250 273 274 275 280 283 288 295

(جدول-۳): درصد خطای بازشناسی واج برای مجموعه توسعه و مجموعه آزمون یا ارزیابی تعریف شده بر روی دادگان فارس‌دات به ازای روش‌های مختلف استخراج ویژگی‌های و مدل‌سازی آکوستیکی در نرم افزار کلدی  
(Table-3): Percent of Phoneme error rate (PER) on the development & evaluation datasets defined on the FarsDat database for applying various methods of feature extraction and acoustic modeling in the Kaldi Toolkit

شماره	ویژگی	مدل	خطای بازشناسی واج (%)	
			مجموعه توسعه	مجموعه آزمون
1	MFCC + delta + delta-delta	Mono-phone	29.8	30.5
2	MFCC + delta + delta-delta	Tri-phone	25.8	25.1
3	MFCC + LDA + MLLT	Tri-phone	24.4	23.3
4	MFCC + LDA + MLLT + SAT	Tri-phone	21.2	21.8
5	MFCC + LDA + MLLT + SAT	SGMM	20.2	19.9
6	MFCC + LDA + MLLT + SAT	SGMM + MMI	20.3	19.8
7	MFCC + LDA + MLLT + SAT	Karel-DNN	20.3	19.8

## ۵- جمع بندی

در این مقاله، نگارنده با اعتقاد به این که در حال حاضر سامانه متن‌باز بازشناسی گفتار کلدی، مدرن‌ترین و کارآمدترین سامانه حال حاضر دنیاست، ابتدا به بررسی ویژگی‌ها و توانمندی‌های مختلف آن پرداخته است. یکی از ویژگی‌های منحصر به فرد کلدی، وجود کدهای آماده در آن برای انجام بازشناسی بر روی دادگان‌های گفتاری متداول به زبان‌های مختلف است؛ که کمک بزرگی به پژوهش‌گران در بازشناسی نتایج قبلی و ارائه نتایج جدید در یک بستر یکسان و قابل دسترس برای همه محسوب می‌شود. در راستای بنای یک بستر همگانی و استاندارد برای پژوهش‌گران فعال در عرصه بازشناسی گفتار فارسی، در این مقاله بر مبنای دادگان فارس‌دات و با الهام از کار صورت‌گرفته برای دادگان TIMIT، مجموعه‌های آموزش، توسعه و ارزیابی انتخاب شده و دقت بازشناسی با اعمال تکنیک‌های مختلف موجود در کلدی محاسبه شده است. بهترین میزان خطای حاصل در بازشناسی واج برای مجموعه توسعه ۲۰/۳ درصد و برای مجموعه آزمون ۱۹/۸ بوده است. در این کار، بررسی گسترده اثر به‌کارگیری شبکه‌های عصبی ژرف و سامانه‌های ترکیبی، استخراج ویژگی BNF، به‌کارگیری فرکانس گام برای زبان فارسی در بستر ارائه شده در این مقاله مد نظر است.

## سپاس‌گزاری

نگارنده مقاله بر خود لازم می‌داند که از همکاری و حمایت بی‌دریغ آقای دکتر دنیل پووی (Daniel Povey) پژوهش‌گر ارشد مرکز پردازش زبان و گفتار (CLSP) دانشگاه جانز هاپکینز، دکتر حسین صامتی عضو هیئت علمی دانشکده مهندسی کامپیوتر دانشگاه صنعتی شریف و مدیر آزمایشگاه پردازش گفتار آن دانشکده، مهندس خسرو حسین‌زاده مدیر فنی شرکت عصرگوش پرداز، دکتر محمد بحرانی عضو هیئت علمی گروه زبان و زبان‌شناسی دانشگاه صنعتی شریف و مدیر بخش پردازش زبان‌های طبیعی شرکت عصرگوش پرداز، دکتر هادی ویسی عضو هیئت علمی دانشکده علوم و فنون دانشگاه تهران و مدیر عامل شرکت عصرگوش پرداز، دکتر دیگو جیولیانی (Diego Giuliani) پژوهش‌گر آزمایشگاه تکنولوژی زبان‌های گفتاری مؤسسه پژوهشی برنوکسلر کشور ایتالیا صمیمانه تشکر کند.

در جدول (۳) نتایج حاصل از انجام بازشناسی مجموعه توسعه و مجموعه آزمون یا ارزیابی به‌ازای ویژگی‌های مختلف استخراج‌شده و روش‌های مختلف مدل‌سازی آکوستیکی آمده است. سطر نخست نتایج را به‌ازای ساده‌ترین حالت یعنی روش استخراج ویژگی MFCC و مشتقات نخست و دوم آن و مدل‌سازی آکوستیکی مبتنی بر تک‌واج نشان می‌دهد. در سطر دوم، فقط مدل‌سازی از تک‌واجی به سه‌واجی تغییر کرده است که طبق انتظار، بهبود درخور توجهی در کاهش خطای بازشناسی واج بر روی هر دو مجموعه توسعه و آزمون را به همراه داشته است. در سطر سوم جدول، در مرحله استخراج ویژگی به جای بردار ویژگی MFCC و مشتقات نخست و دوم آن، بر روی هر بردار ویژگی MFCC سیزده مؤلفه‌ای به همراه سه بردار ویژگی قبل و بعد که تشکیل یک بردار ۹۱ مؤلفه‌ای می‌دهد دو تبدیل خطی LDA و MLLT به ترتیب اعمال شده و در نهایت به یک بردار چهل مؤلفه‌ای می‌رسیم که مدل‌سازی سه‌واجی بر مبنای آن صورت گرفته است. با مقایسه نتایج این مرحله و مرحله قبل، به این نتیجه می‌رسیم که ویژگی MFCC + LDA + MLLT نسبت به MFCC و مشتقات نخست و دوم آن خطای بازشناسی واج را به‌طور مطلق حدود ۱/۵ درصد کاهش داده است و در مراحل بعدی از این ویژگی استفاده خواهیم کرد. در مرحله چهارم هم در مرحله آموزش آن، به کمک تکنیک fMLLR آموزش تطبیقی به گوینده (SAT) صورت گرفته است و هم در مرحله بازشناسی به روش fMLLR، ویژگی‌های هر گوینده تطبیق داده شده است که سرجمع بر روی مجموعه توسعه ۳/۲ درصد و بر روی مجموعه آزمون ۱/۵ درصد به‌طور مطلق بهبود دقت به همراه داشته است. در مرحله پنجم به‌کارگیری مدل‌سازی SGMM موجب بهبود یک‌درصدی بر روی مجموعه توسعه و بهبود حدود دودرصدی بر روی مجموعه آزمون شده است. و در نهایت اعمال آموزش تمایزی<sup>۱</sup> MMI [26] حدود یک‌صدم درصد بهبود به همراه داشته است.

در نهایت در سطر آخر جدول نتیجه به‌کارگیری شبکه عصبی ژرف برای مدل‌سازی آکوستیکی سه‌واجی‌ها آمده که نسبت به بهترین دقت حاصله در دو مرحله قبل هیچ‌گونه بهبودی به همراه ندارد که دلیل آن هم کوچک بودن اندازه مجموعه آموزشی است. در کارهای آتی به تبیین اثر به‌کارگیری شبکه‌های عصبی ژرف در بازشناسی گفتار فارسی بر روی دادگان‌های بزرگتر فارسی خواهیم پرداخت.

<sup>۱</sup> Maximum Mutual Information

- [10] F. Seide, G. Li, and D. Yu, Conversational speech transcription using context-dependent deep neural networks. in Proc. of Interspeech, pp. 437-440, 2011.
- [11] K. Vesely, A. Ghoshal, L. Burget, and D. Povey, Sequence discriminative training of deep neural networks, in Proc. of Interspeech, pp. 2345-2349, 2013.
- [12] D. Povey, X. Zhang, and S. Khudanpur, Parallel training of Deep Neural Networks with Natural Gradient and Parameter Averaging, in Proc. of 3rd International Conference on Learning Representations (ICLR2015), USA, 2015
- [13] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, Recurrent neural network based language model, in Proc. of Interspeech, pp. 1045-1048, Japan, 2010
- [14] S. Young, G. Evermann, M. Gales, T. Hain, D. Kershaw, X. Liu, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, and P. Woodland, The HTK Book (for version 3.4). Cambridge Univ. Eng. Dept., 2009.
- [15] B. Pellom, SONIC: The University of Colorado Continuous Speech Recognizer, Technical Report TRCSLR-2001-01, Center for Spoken Language Research, University of Colorado, USA, 2001.
- [16] B. Pellom and K. Hacioglu, Recent Improvements in the CU SONIC ASR System for Noisy Speech: The SPINE Task, in Proc. of ICASSP, Hong Kong, 2003.
- [17] K.F. Lee, H.W. Hon, and R. Reddy, An overview of the SPHINX speech recognition system, in IEEE Transactions on Acoustics, Speech and Signal Processing 38.1, pp.35-45, 1990.
- [18] W. Walker, P. Lamere, P. Kwok, B. Raj, R. Singh, E. Gouvea, P. Wolf, and J. Woelfel, Sphinx-4: A flexible Open Source Framework for Speech Recognition, Sun Microsystems Inc., Technical Report SML1 TR2004-0811, 2004.
- [19] D. Rybach, C. Gollan, G. Heigold, B. Hoffmeister, J. Löff, R. Schlüter, and H. Ney, The RWTH Aachen University Open Source Speech Recognition System, in Proc. of Interspeech, pp. 2111-2114, 2009.
- [20] D. Rybach, S. Hahn, P. Lehnen, D. Nolden, M. Sundermeyer, Z. Tüske, S. Wiesler, R. Schlüter, and H. Ney: RASR - The RWTH Aachen University Open Source Speech Recognition Toolkit, In IEEE Automatic Speech Recognition

## 6-Refrence

۶-مراجع

- [1] B. BabaAli and H. Sameti, The Sharif speaker-independent large vocabulary speech recognition system, in Proceedings of the 2nd Workshop on Information Technology & Its Disciplines (WITID '04), pp. 24-26, Iran, 2004.
- [2] H. Sameti, H. Veisi, M. Bahrani, B. Babaali, K. Hosseinzadeh, A large vocabulary continuous speech recognition system for Persian language, in EURASIP Journal on Audio, Speech, and Music Processing 2011, 2011:6
- [3] F. Almasganj, S.A. Seyyed Salehi, M. Bijankhan, H. Razizade, M. Asghari, Shenava 2: a Persian continuous speech recognition software, in The first workshop on Persian language and Computer, pp. 77-82, Tehran, 2004.
- [4] M. Sheikhan, M. Tebyani, M. Lotfizad, Continuous speech recognition and syntactic processing in iranian farsi language. Inter J Speech Technol. 1(2), 135 (1997). doi:10.1007/BF02277194
- [5] S.M. Ahadi, Recognition of continuous Persian speech using a medium-sized vocabulary speech corpus, in European Conference on Speech communication and technology (Eurospeech'99), pp. 863-866, Switzerland, 1999.
- [6] N. Srinivasamurthy, S.S. Narayanan, Language-adaptive Persian speech recognition, in European Conference on Speech Communication and Technology (Eurospeech'03), Switzerland, 2003.
- [7] H. Sameti, H. Veisi, M. Bahrani, B. Babaali, K. Hosseinzadeh, Nevisa, a Persian continuous speech recognition system, in Communications in Computer and Information Science (Springer Berlin Heidelberg), pp. 485-492, 2008.
- [8] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, J. Silovsky, G. Stemmer, and K. Vesely, The Kaldi Speech Recognition Toolkit, in IEEE 2011 Workshop on Automatic Speech Recognition and Understanding. IEEE Signal Processing Society, 2011.
- [9] G. Hinton, L. Deng, D. Yu, G.E. Dahl, A.-R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. Sainath et al., Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups, Signal Processing Magazine, IEEE, vol. 29, no. 6, pp. 82-97, 2012.

- for LVCSR of meetings, in Proc. ICASSP, pp. 757–760, 2007.
- [32] M. J. F. Gales, Maximum likelihood linear transformations for HMM based speech recognition, in Computer Speech and Language, vol. 12, no. 2, pp. 75–98, 1998.
- [33] S. Furui, Cepstral analysis technique for automatic speaker verification, in IEEE Transactions on Acoustics, Speech, and Signal Processing, vol. 29, no. 2, pp. 254–272, 1981.
- [34] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, Front-end factor analysis for speaker verification, in IEEE Transactions on Audio, Speech, and Language Processing, vol. 19, no. 4, pp. 788–798, 2011.
- [35] L. Lee and R. Rose, Speaker Normalization Using Efficient Frequency Warping Procedure, in Proc. of ICASSP, pp. 1–353–356, Atlanta, USA, 1996,
- [36] D. Kim, S. Umesh, M. Gales, T. Hain, and P. Woodland, Using VTLN for Broadcast News Transcription, In Proc. of the 8th ICSLP, Jeju Island, Korea, 2004.
- [37] L. Burget, Combination of Speech Features Using Smoothed Heteroscedastic Linear Discriminant Analysis. In Proc. of the 8th ICSLP, Jeju Island, Korea, pp. 2549–2552, 2004
- [38] K. Visweswariah, S. Axelrod, R.A. Gopinath, Acoustic modeling with mixtures of subspace constrained exponential models. In Proc. of the 7th Eurospeech'2003, Geneva, Switzerland, pp. 2613–2616, 2003.
- [39] Xin Lei, Modeling Lexical Tones for Mandarin Large Vocabulary Continuous Speech Recognition, Ph.D. thesis, University of Washington, 2006.
- [40] P. Ghahremani, B. BabaAli, D. Povey, K. Riedhammer, J. Trmal, and S. Khudanpur, A pitch extraction algorithm tuned for automatic speech recognition, In Proc. of ICASSP, pp. 2494–2498, Italy, 2014.
- [41] S. Renals, N. Morgan, H. Bourlard, M. Cohen, and H. Franco, Connectionist probability estimators in HMM speech recognition, in IEEE Trans. Speech Audio Process., vol. 2, no. 1, pp. 161–174, 1994.
- [42] A. J. Robinson, An application of recurrent nets to phone probability estimation,” in IEEE Trans. on Neural Networks, vol. 5, no. 2, pp. 298–305, 1994.
- and Understanding Workshop (ASRU), USA, 2011.
- [21] A. Lee, T. Kawahara, and K. Shikano, JULIUS - an open source real-time large vocabulary recognition engine, in Proc. of INTERSPEECH, pp. 1691–1694, 2001.
- [22] K. Demuynck, J. Roelens, D.V. Compernelle and P. Wambacq, SPRAAK: An Open Source SPeech Recognition and Automatic Annotation Kit, In Proc. Interspeech, pp. 495–498, Australia, 2008.
- [23] D. Bolaños, The Bavieca Open-Source Speech Recognition Toolkit, in Proc. of IEEE Workshop on Spoken Language Technology (SLT), Miami, FL, USA, 2012.
- [24] M. Bijankhan and M. J. Sheikhzadegan, FARSDAT—the speech database of Farsi spoken language, in Proceedings of the 5th Australian International Conference on Speech Science and Technology (SST '94), pp. 826–829, Perth, Australia, December 1994.
- [25] V. Zue, S. Seneff, and J. Glass, Speech database development at MIT: TIMIT and beyond, Speech Communication, vol. 9, no. 4, pp. 351–356, 1990.
- [26] K.F. Lee, and H.W. Hon, Speaker-independent phone recognition using hidden Markov models, in IEEE Transactions on Acoustics, Speech and Signal Processing, vol 37, no. 11, pp. 1641–1648, 1989.
- [27] M. Mohri, F. Pereira, and M. Riley, Weighted finite-state transducers in speech recognition, in Computer Speech & Language 16, no. 1, pp. 69–88, 2002.
- [28] C. Allauzen, M. Riley, J. Schalkwyk, W. Skut, and M. Mohri, OpenFst: a general and efficient weighted finite-state transducer library, in Proc. CIAA, 2007.
- [29] D. Povey, L. Burget, M. Agarwal, P. Akyazi, F. Kai, A. Ghoshal, O. Glembek, N. Goel, M. Karafit, A. Rastrow, R. C. Rose, P. Schwarz, and S. Thomas, The subspace Gaussian mixture model – A structured model for speech recognition, in Computer Speech & Language, vol. 25, no. 2, pp. 404 – 439, 2011.
- [30] H. Hermansky, Perceptual linear predictive (PLP) analysis of speech, Journal of the Acoustical Society of America, vol. 87, no.4, pp. 1738–1752, 1990.
- [31] F. Grezl, M. Karafiat, S. Kontar, and J. Cernocky, Probabilistic and bottle-neck features



**باقر باباعلی** مدرک کارشناسی خود را در سال ۱۳۷۹ در رشته مهندسی کامپیوتر (گرایش سخت‌افزار) از دانشگاه شیراز اخذ کرد. همچنین مدرک کارشناسی ارشد و دکترای خود را به ترتیب در سال‌های ۱۳۸۲ و ۱۳۸۹ در رشته مهندسی کامپیوتر (گرایش هوش مصنوعی) از دانشگاه صنعتی شریف دریافت کرد. زمینه‌های پژوهشی مورد علاقه ایشان یادگیری ماشین و بازشناسی الگوی آماری، بازشناسی گفتار، بازشناسی الگوهای دنباله‌ای، یادگیری ژرف و کاربردهای آن بوده و در حال حاضر عضو هیأت علمی در دانشکده ریاضی، آمار و علوم کامپیوتر دانشگاه تهران است. نشانی رایانامه ایشان عبارت است از:

**babaali@ut.ac.ir**

- [43] G. E. Dahl, D. Yu, L. Deng, and A. Acero, Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition, in *IEEE Trans. Audio Speech Lang. Process.*, vol. 20, no. 1, pp. 30–42, 2012.
- [44] N. Jaitly, P. Nguyen, A. Senior, and V. Vanhoucke, Application of pretrained deep neural networks to large vocabulary speech recognition, in *Proc. INTERSPEECH*, September 2012.
- [45] T. N. Sainath, B. Kingsbury, B. Ramabhadran, P. Fousek, P. Novak, and A. Mohamed, Making deep belief networks effective for large vocabulary continuous speech recognition, in *Proc. IEEE ASRU*, pp. 30–35, December 2011.
- [46] G. E. Hinton, S. Osindero, and Y. Teh, A fast learning algorithm for deep belief nets, *Neural Computation*, vol. 18, pp. 1527–1554, 2006.
- [47] M. Gales, The generation and use of regression class trees for MLLR adaptation, University of Cambridge, Department of Engineering, 1996.
- [48] D. Povey, G. Zweig, and A. Acero, The Exponential Transform as a generic substitute for VTLN, in *Proc. IEEE ASRU*, December 2011.
- [49] S.J. Young, P.C. Woodland, The use of state tying in continuous speech recognition, in *European Conference on Speech Communication and Technology (EUROSPEECH'93)*, pp. 2203–2206, Germany, September 1993.
- [50] M. Federico, N. Bertoldi, and M. Cettolo, IRSTLM: an Open Source Toolkit for Handling Large Scale Language Models, in *Proc. of Interspeech*, Brisbane, Australia, 2008.
- [51] A. Stolcke, SRILM-an extensible language modeling toolkit, in *Proc. of INTERSPEECH*, 2002.
- [52] D. Povey, D. Kanevsky, B. Kingsbury, B. Ramabhadran, G. Saon, and K. Visweswariah, Boosted MMI for model and feature-space discriminative training, in *Proc. of ICASSP*, pp. 4057–4060, 2008.