



# بازشناسی دست‌نوشته برون خط عربی بر مبنای یک رویکرد تلفیقی جدید از مدل مخفی مارکوف و شبکه‌های عصبی ژرف

باقر باباعلی\* و بابک رکابدار

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

## چکیده

مسئله مدل‌سازی و بازشناسی دست‌نوشته شباهت بسیار زیادی به مسئله مدل‌سازی و بازشناسی گفتار دارد. به همین علت می‌توان از رویکردهای به کار گرفته شده برای مسئله بازشناسی گفتار با اندکی تغییر در مراحل ابتدایی آن مانند استخراج ویژگی، برای بازشناسی دست‌نوشته نیز بهره برد. با گسترش رویکردهای ترکیبی HMM-DNN و استفاده از توابع هدف دنباله‌ای مانند MMI پیشرفت‌های قابل توجهی در حوزه بازشناسی گفتار حاصل شده است. این مقاله با استفاده از نرم‌افزار متن‌باز KALDI، که شهرت اصلی آن در حوزه بازشناسی گفتار و همچنین به کارگیری آخرین مدل‌های ترکیبی ارائه شده در آن، به کمک روش افزایش داده مدلی برای بازشناسی دست‌نوشته عربی ارائه داده است. این پژوهش بر روی دادگان KHA'TT انجام شده که نرخ خطای بازشناسی واژه را بر روی این دادگان به میزان ۷/۳۲ درصد مطلق کاهش داده است.

واژگان کلیدی: بازشناسی دست‌نوشته عربی، شبکه‌های عصبی ژرف، مدل مخفی مارکوف، نرم‌افزار متن‌باز KALDI

## Off-line Arabic Handwritten Recognition Using a Novel Hybrid HMM-DNN Model

Bagher BabaAli\* & Babak Rekadardar

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

### Abstract

In order to facilitate the entry of data into the computer and its digitalization, automatic recognition of printed texts and manuscripts is one of the considerable aid to many applications. Research on automatic document recognition started decades ago with the recognition of isolated digits and letters, and today, due to advancements in machine learning methods, efforts are being made to identify a sequence of handwritten words. Generally, based on the type of text, document recognition is divided into two main categories: printed and handwritten. Due to the limited number of fonts relative to the diversity of handwriting of different writers, it is much easier to recognize printed texts than handwritten text; thus, the technology of recognizing printed texts has matured and has been marketed in the form of a product. Handwriting recognition task is usually done in two ways: online and offline; offline handwriting recognition involves the automated translation of text in image format to letters that can be used in computer and text-processing applications. Most of the research in the field of handwriting recognition has been conducted on Latin script, and a variety of tools and resources have been gathered for this script. This article focuses on the application of the latest methods in the field of speech recognition for the recognition of Arabic handwriting. The task of handwritten text modeling and recognizing is very similar to the task of speech modeling and recognition. For this reason, it is possible to apply the approaches used for the speech recognition with a slight change for the handwriting recognition. With the expansion of HMM-DNN hybrid approaches and

\* Corresponding author

\*نویسنده عهده‌دار مکاتبات

the use of sequential objective functions such as MMI, significant improvements have been made in the accuracy of speech recognition system. This paper presents a pipeline for the offline Arabic handwritten text recognition using the open source KALDI toolkit, which is very well-known in the community of speech recognition, as well as the use of the latest hybrid models presented in it and data augmentation techniques. This research has been conducted on the Arabic KHATT database, which achieved 7.32% absolute reduction in word recognition error (WER) rate.

**Keywords:** Arabic Handwritten Recognition, Deep Neural Networks, Hidden Markov Model, Kaldi Toolkit

باشد که در این صورت، به طور معمول لازم است در ابتدا قطعه‌بندی خطوط صورت گرفته و در ادامه بازنشاسی تک تک خطوط به طور مجزا صورت گیرد. کیفیت تصاویر حاصل از پویش دست‌نوشته‌ها، تنوع بسیار بالای دست‌نوشته افراد مختلف، فاصله‌های ایجاد شده بین کلمات و یا حتی بخش‌هایی از یک کلمه و همچنین تنوع در اندازه و نوع قلم را می‌توان از جمله عواملی دانست که بازنشاسی دست‌نوشته برون خط را به یکی از چالش‌برانگیزترین پژوهش‌ها در حوزه تحلیل و بازنشاسی متون تبدیل کرده است.



(شکل-۱): بازنشاسی دست‌نوشته به دو صورت برخط و برون خط  
(Figure-1): Offline and online handwritten recognition

خطوط نگارشی زنده و متداول در دنیا شامل لاتین، سیریلیک، عربی، هندی، گرجی و ارمنی است که متداول‌ترین آن لاتین است و از همین رو عمده پژوهش‌ها بیشتر بر روی خط لاتین صورت گرفته و منابع و دادگان‌های متعددی در این زمینه گردآوری شده است. در این بین، پژوهش‌های به مراتب کمتری بر روی دست‌نوشته عربی و به طور کلی زبان‌هایی که از خط عربی بهره می‌برند، از جمله فارسی رایج در ایران انجام شده است. این پژوهش با هدف پرداختن به بازنشاسی دنباله‌ای از کلمات دست‌نوشته پیوسته فارسی بر مبنای رویکردها و روش‌های نوین ارائه شده در سالیان اخیر آغاز شد، ولی به دلیل عدم وجود یک دادگان دست‌نوشته فارسی مناسب برای انجام آزمایش‌ها و همچنین قرابت زیاد خط فارسی به عربی، آزمایش‌ها بر مبنای خط عربی و به طور مشخص دادگان KHATT انجام شد.

<sup>6</sup> Line Segmentation

## ۱- مقدمه

در راستای تسهیل ورود اطلاعات به رایانه و دیجیتال‌سازی آن، بازنشاسی خودکار متون چاپی و دست‌نوشته کمک بزرگی محسوب می‌شود که کاربردهای متعددی نظیر دسته‌بندی خودکار نامه‌ها، خواندن خودکار چک‌های بانکی و فرم‌های پر شده دستی، تبدیل متون دست‌نوشته تاریخی و صدها کاربرد دیگر برای آن قابل تصور است. پژوهش‌ها بر روی بازنشاسی ماشینی متون از چند دهه قبل، ابتدا با بازنشاسی ارقام و حروف محدود و مجزا آغاز شد و امروزه به لطف پیشرفت‌های علمی حاصل شده در حوزه یادگیری ماشین و بازنشاسی الگو به تلاش در جهت بازنشاسی دنباله‌ای از کلمات دست‌نوشته پیوسته منجر شده است. بازنشاسی متون بر مبنای نوع متن به دو صورت چاپی<sup>۱</sup> و دست‌نوشته<sup>۲</sup> تقسیم‌بندی می‌شود. بنا به محدود بودن تعداد فونت‌های چاپی در مقایسه با تنوع دست‌خط افراد مختلف، بازنشاسی متون چاپی در مقایسه با دست‌نوشته به مراتب راحت‌تر است؛ به نحوی که فناوری تشخیص متون چاپی به بلوغ رسیده و در قالب محصول به بازار عرضه شده است. تشخیص دست‌نوشته به طور معمول به دو صورت برخط<sup>۳</sup> و برون خط<sup>۴</sup> انجام می‌شود. در حالت برخط، کاربر بر روی یک لوح فشرده شروع به نوشتن می‌کند و مختصات نوک قلم با سرعت ثابتی (به طور معمول یکصد تا دویست هرتز) نمونه‌برداری<sup>۵</sup> شده و بر مبنای این دنباله از مختصات نوک قلم، بازنشاسی دست‌نوشته صورت می‌گیرد؛ ولی در حالت برون خط، ورودی تصویری است که از دست‌نوشته گرفته شده و از اطلاعات پویا حین نوشتن بی‌بهره است. به قدری این اطلاعات پویا در بهبود دقت بازنشاسی مؤثر است که در [۱۱] تلاش در جهت بازسازی آن از روی تصویر صورت گرفته است. در حالت برون خط تصویر دست‌نوشته ممکن است شامل خطوط متعددی (به عنوان مثال یک پاراگراف)

<sup>1</sup> Printed

<sup>2</sup> Handwritten

<sup>3</sup> Online

<sup>4</sup> Offline

<sup>5</sup> Sampling

در ادامه این مقاله، در بخش بعد به مروری اجمالی بر روی دادگان‌ها و پژوهش‌های انجام‌شده در حوزه بازشناسی دست‌نوشته فارسی و عربی می‌پردازیم. بخش سوم به بررسی نحوه مدل‌سازی بازشناسی دست‌نوشته و رویکردهای رایج آن اختصاص داده شده است. در بخش چهارم به تشریح اجزای سامانه به‌کار گرفته‌شده در این پژوهش پرداخته‌ایم. نتایج آزمایش‌ها و بررسی و تحلیل آنها در بخش پنجم آمده و در نهایت در بخش ششم جمع‌بندی ارائه شده است.

## ۲- مرور اجمالی دادگان‌ها و پژوهش‌ها

غالب پژوهش‌های انجام‌شده در حوزه بازشناسی دست‌نوشته فارسی به بازشناسی ارقام و حروف مجزا محدود می‌شود. در [34] و [31] و [36] با استفاده از روش‌های مختلف استخراج ویژگی و همچنین بهره‌گیری از روش‌های دسته‌بندی ماشین بردار پشتیبان<sup>۱</sup> و  $K-NN^2$  تلاش‌هایی برای بازشناسی ارقام دست‌نوشته فارسی انجام و در [24] همین هدف با استفاده از تجزیه ساختاری<sup>۳</sup> دنبال شده است؛ همچنین پژوهش انجام‌شده در [18] برای بازشناسی ارقام دست‌نوشته فارسی از روش‌های بینایی رایانه<sup>۴</sup> و ویژگی‌های توصیفی تصویر مانند  $HOG^5$  بهره برده است. اما در کنار این پژوهش‌ها، موارد محدودی نیز برای بازشناسی کلمات مجزا در زبان فارسی نیز گزارش شده است. [8] و [9] تلاش‌هایی جهت بازشناسی کلمه توسط مدل مخفی مارکوف بوده‌اند؛ همچنین [6] با استخراج بردار ویژگی از تصاویر کلمات فارسی توسط تبدیل ویولت به‌دنبال بازشناسی کلمات فارسی از روی تصاویر بوده است. در [43] یک روش دومرحله‌ای برای بازشناسی کلمات دست‌نوشته فارسی به‌کمک بلوک‌بندی تطبیقی گرادینان تصویر ارائه شده است. در همین‌اواخر نیز [30] از شبکه‌های عصبی ژرف<sup>۶</sup> برای بازشناسی کلمات دست‌نوشته مجزا استفاده کرده است.

تا جایی که مطلع‌ایم، هیچ پژوهشی روی بازشناسی دنباله‌ای از کلمات فارسی به‌صورت یک خط دست‌نوشته گزارش نشده است. همان‌طور که در قبل گفته شد، مهم‌ترین علت این امر نبود دادگان دست‌نوشته فارسی مناسب برای این منظور بوده است. در زبان عربی نیز تعداد دادگان برای

این منظور بسیار محدود است. کمبود دادگان مناسب باعث محدودیت پژوهش در این حوزه در زبان عربی و فارسی شده است. برای رفع این معضل، تلاش‌هایی به‌منظور جمع‌آوری دادگان استاندارد برای بازشناسی دست‌نوشته فارسی و عربی صورت گرفته است. هر چند در این پژوهش هدف بررسی بازشناسی دنباله‌ای از کلمات دست‌نوشته است، ولی در ادامه به بررسی مهم‌ترین دادگان‌های موجود عربی و فارسی که در بردارنده حروف مجزا، ارقام مجزا و رشته‌ای، تاریخ، واژه مجزا و در نهایت دنباله کلمات (متن) هستند، می‌پردازیم.

یکی از آخرین تلاش‌ها برای جمع‌آوری دادگان فارسی توسط [32] انجام شده است. این دادگان از تصاویر ارقام، حروف، واژگان و متون دست‌نوشته فارسی تشکیل شده است. تعداد نویسندگان این دادگان پانصد نفر است که از ۲۵۰ نفر مرد و ۲۵۰ نفر زن تشکیل شده است. مشکل عمده این دادگان نبود قطعه‌بندی خطی برای متون نوشته‌شده در قالب پاراگراف و همچنین نبود متن متناظر با هر خط از دست‌نوشته است که باعث شده قابل به‌کارگیری در پژوهش جاری نباشد.

یکی از مشهورترین دادگان دست‌نوشته فارسی دادگان HODA [19] است. دادگان HODA شامل ارقام دست‌نویس صفر تا نه است و از تصاویر دست‌نوشته بیش از ۱۱۹۴۲ نفر تشکیل شده است. تعداد تصاویر این دادگان بالغ بر هشتاد هزار تصویر است که شصت هزار تعداد از آنها به‌عنوان مجموعه آموزش و بیست هزار تعداد از آنها به‌عنوان مجموعه آزمون در نظر گرفته شده است. یکی دیگر از مهمترین دادگان‌های فارسی، دادگان IFN/Farsi [22] است. این دادگان از تصاویر اسامی شهرها و روستاهای ایران تشکیل شده و برای مواردی که هدف بازشناسی کلمات مجزای فارسی است، مناسب است. تعداد نویسندگان آن ششصد نفر و تعداد تصاویر آن ۷۲۷۱ تصویر است.

تنها دادگان فارسی که شامل تصاویر خطوط دست‌نوشته فارسی است، و یا به عبارتی از قطعه‌بندی خطی برخوردار است دادگان HaFT [33] است. این دادگان با همکاری ششصد نویسنده تهیه شده است. هر نویسنده، سه برگه را که هر کدام شامل هشت خط بوده، نوشته است. در این دادگان تصاویر هر خط به‌صورت جداگانه در دسترس است؛ اما به‌دلیل این که متن هر تصویر موجود نیست مشابه دادگان [32] نمی‌توان از آن برای بازشناسی دنباله کلمات دست‌نوشته پیوسته فارسی که هدف این مقاله است، بهره برد.

<sup>1</sup> Support Vector Machine

<sup>2</sup> K-Nearest neighbor

<sup>3</sup> Structural Decomposition

<sup>4</sup> Computer Vision

<sup>5</sup> Histogram of Oriented Gradients

<sup>6</sup> Deep Neural Network

از مرسوم‌ترین دادگان‌های عربی می‌توان به IFN/Enit [25] و KHATT [21] اشاره کرد که اولی شامل کلمات مجزا و دومی تصاویر خطوط دسته‌نوشته است که هر خط شامل دنباله‌ای از واژگان دست‌نوشته است. دادگان IFN/Enit شامل ۲۶۴۵۹ تصویر دست‌نوشته از ۹۳۷ اسم شهر و روستا در کشور تونس تهیه شده است. تعداد نویسندگان این دادگان ۴۱۱ نفر است. دادگان KHATT شامل ۱۳۳۵۷ تصویر از خطوط دست‌نوشته این دادگان شامل دوهزار پاراگراف یکسان و دوهزار پاراگراف منحصر به فرد است. این ویژگی باعث می‌شود که بتوان از این

دادگان برای مسائل دیگری مانند تشخیص نویسنده نیز بهره برد. تعداد هزار نویسنده برای جمع‌آوری این دادگان مشارکت داشته‌اند. دادگان KHATT تنها دادگان عربی در دسترس است که برخلاف دو دادگان فارسی موجود، هم قطعه‌بندی خطی و هم متن متناظر با هر خط را دارد، به همین دلیل این دادگان به‌عنوان تنها گزینه موجود برای انجام آزمایش‌های این پژوهش انتخاب شد. تعداد بیشتری از دادگان‌های حوزه بازشناسی دست‌نوشته فارسی و عربی در جدول (۱) ادامه است.

(جدول-۱): دادگان‌های متداول برای بازشناسی دست‌نوشته فارسی و عربی

(Table-1): Popular datasets for Persian/Arabic handwritten recognition

تعداد تصاویر	تعداد نویسنده	نوع داده	زبان	نام دادگان
250624	500	ارقام، حروف، کلمات و متون دست‌نوشته	فارسی	Sadri2016 [32]
80000	11942	ارقام دست‌نوشته	فارسی	HODA [19]
7271	600	اسامی شهرها و روستاهای ایران	فارسی	IFN/Farsi [22]
14400	600	یک خط دست‌نوشته	فارسی	IlaFT [33]
1000	250	پاراگراف دست‌نوشته	فارسی	FHT [42]
34200	380	اسامی شهرهای ایران	فارسی	IAUT/PLICN [5]
70120	-	حروف و ارقام دست‌نوشته	فارسی	IFHCDB [23]
41004	295	ارقام، حروف، کلمات و متون دست‌نوشته	فارسی	CENPARMI [15]
13357	1000	یک خط دست‌نوشته	عربی	KHATT [21]
26459	411	اسامی شهرها و روستاهای تونس	عربی	IFN/Enit [25]
750000	305	یک خط دست‌نوشته	عربی	MADCAT Arabic [20]

### ۳- مدل‌سازی بازشناسی دست‌نوشته

مسئله مدل‌سازی و بازشناسی دست‌نوشته شباهت بسیار زیادی به مسئله بازشناسی گفتار دارد. از همین رو سامانه‌های متن‌باز موجود در حوزه بازشناسی گفتار نظیر KALDI [26]، HTK [41] و RASR [29] به‌منظور بازشناسی دست‌نوشته نیز به‌کار گرفته شده‌اند. اغلب ابتدا روش‌ها برای مسئله بازشناسی گفتار توسعه داده می‌شوند و سپس برای حوزه بازشناسی دست‌نوشته اختصاصی می‌شوند. جدول (۲) اجزای و مسائل مشابه این دو حوزه بازشناسی را در هر سطر جدول ارائه داده است.

(جدول-۲): اجزای و مسائل مشابه در بازشناسی گفتار و

بازشناسی دست‌نوشته

(Table-2): Similar issues in speech recognition and handwritten recognition

حوزه بازشناسی گفتار	حوزه بازشناسی دست‌نوشته
مدل‌سازی آکوستیک	مدل‌سازی ویژوال

مدل‌سازی زبانی	مدل‌سازی زبانی
مدل‌سازی طول واج	مدل‌سازی طول شناسه
مدل‌سازی بافت واج	مدل‌سازی بافت شناسه
تطبیق گوینده	تطبیق نویسنده

مشابه بازشناسی گفتار، در رویکرد آماری بازشناسی دست‌نوشته، با بیشینه‌سازی احتمال پسین رابطه (۱) و بهره‌گیری از قاعده بیز محتمل‌ترین دنباله واژگان (متن) به‌ازای تصویر دست‌نوشته ورودی به‌دست می‌آید:

$$W^* = \underset{W}{\operatorname{argmax}} P(W|X) = \underset{W}{\operatorname{argmax}} P(X|W)P(W) \quad (1)$$

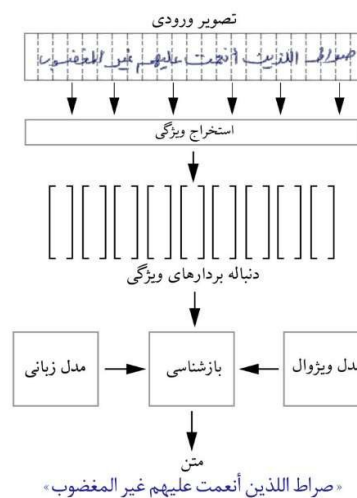
که  $X = x_1 x_2 x_3 \dots x_T$  دنباله بردارهای ویژگی استخراج‌شده از دست‌نوشته و  $W = w_1 w_2 \dots w_T$  دنباله واژگان متن است. در این رابطه  $P(X|W)$  مدل ویژوال است که توزیع احتمالی دست‌نوشته به شرط متن را تخمین می‌زند و

بردارهای ویژگی استخراج‌شده بر مبنای پیکسل از تصاویر دودویی دست‌نوشته و همچنین استفاده از ایده ترکیب<sup>۶</sup> نتایج خروجی سه سامانه بازشناسی مبتنی بر HMM سامانه کارایی ارائه داده شده است. در [3] که هدف بازشناسی کلمات مجزا است دو نوع ویژگی از تصویر واژه استخراج می‌کند. ویژگی نخست بر مبنای پنجره‌ای که بر روی تصویر لغزانیده داده می‌شود، استخراج می‌شود و از آن برای آموزش مدل‌های HMM استفاده می‌کند. ویژگی دوم شامل ویژگی‌های ساختاری مانند تعداد زیرواژه و تعداد علامت‌ها از جمله نقطه است. در این سامانه ابتدا بازشناسی به‌وسیله مدل‌های HMM انجام و در ادامه نتایج به‌دست‌آمده به‌وسیله ویژگی‌های ساختاری امتیازدهی مجدد<sup>۷</sup> می‌شود.

برای سالیان زیادی در عرصه بازشناسی دست‌نوشته روش HMM-GMM به‌عنوان یک رویکرد متداول مطرح بوده است. با مطرح‌شدن یادگیری ژرف در سالیان اخیر، رویکرد معروف به ترکیبی که مبتنی بر تلفیق مدل مخفی مارکوف با شبکه‌های عصبی ژرف (HMM-DNN) است، توانست از رویکرد متداول HMM-GMM پیشی بگیرد و نظر پژوهش‌گران را به خود جلب کند. در رویکرد HMM-GMM برای مدل‌سازی هر شناسه یک مدل مخفی مارکوف چپ به راست آموزش داده می‌شود. برای مدل‌سازی احتمال درست‌نمایی<sup>۸</sup> هر بردار ویژگی از دنباله بردارهای ویژگی استخراج‌شده از دست‌نوشته در هر حالت<sup>۹</sup>، از یک مدل مخلوط گاوسی<sup>۱۰</sup> - که یک مدل تولیدی<sup>۱۱</sup> محسوب می‌شود- بهره گرفته می‌شود. هر مدل مخلوط گاوسی شامل تعدادی گاوسی چندبعدی (هم بعد با بردار ویژگی) است که اوزان، میانگین و ماتریس کوواریانس (به‌طورمعمول قطری در نظر گرفته می‌شود) آنها حین فرآیند آموزش بر مبنای داده‌های آموزشی مربوطه تخمین زده می‌شود.

در رویکرد ترکیبی (HMM-DNN) نحوه به‌کارگیری مدل مخفی مارکوف مشابه رویکرد HMM-GMM است با این تفاوت که برای تخمین احتمال درست‌نمایی بردارهای ویژگی در حالت‌ها، یک شبکه عصبی - که یک روش تمایزی<sup>۱۱</sup> محسوب می‌شود- به‌کار گرفته می‌شود. ورودی این شبکه یک بردار ویژگی به‌همراه تعدادی از بردارهای ویژگی مجاور آن (از هر دو طرف) و خروجی آن احتمال

$P(W)$  مدل زبانی است که احتمال رخداد دنباله واژگان را تخمین می‌زند. در شکل (۲) نمای کلی یک سامانه بازشناسی دست‌نوشته پیوسته آمده است.



(شکل-۲): نمای کلی از سامانه بازشناسی دست‌نوشته  
(Figure-2): Generic block diagram of the handwritten recognition systems

در مدل‌سازی ویژوال  $P(X|W)$ ، هدف برچسب‌زنی دنباله بردارهای ویژگی به طول  $T$  با یک دنباله به طول  $N$  از شناسه‌ها یا واژگان است که به‌طورمعمول  $N$  بسیار کوچک‌تر از  $T$  است. مدل مخفی مارکوف به‌دلیل سادگی و توان بالایی که در مدل‌سازی الگوهای متوالی<sup>۱</sup> دارد، یکی از متداول‌ترین روش‌های مدل‌سازی ویژوال محسوب می‌شود. در این روش، به‌طورمعمول از مدل مخلوط گاوسی<sup>۲</sup> برای مدل‌سازی احتمال مشاهده بردارهای ویژگی در حالات مدل مخفی مارکوف استفاده می‌شود که در اصطلاح رویکرد HMM-GMM نامیده می‌شود.

از مهم‌ترین پژوهش‌هایی که با استفاده از رویکرد HMM-GMM انجام شده است، می‌توان به [16] اشاره کرد. در این پژوهش ابتدا کلمات متن با تجزیه مورفولوژیکی<sup>۳</sup> به زیرکلمات<sup>۴</sup> شکسته می‌شوند؛ بنابراین مجموعه‌واژگان<sup>۵</sup> این سامانه ترکیبی از کلمات و زیرکلمات است که این سبب می‌شود تا با ترکیب زیرکلمات موجود در مجموعه‌واژگان بتوان هر کلمه‌ای را شناسایی کرد. این پژوهش از یک مدل زبانی  $n$ -gram که بر روی متن تجزیه‌شده آموزش داده شده است استفاده می‌کند. همچنین در [2] با استفاده از

<sup>6</sup> Fusion

<sup>7</sup> Rescoring

<sup>8</sup> Likelihood

<sup>9</sup> State

<sup>10</sup> Generative

<sup>11</sup> Discriminative

<sup>1</sup> Sequential Patterns

<sup>2</sup> Gaussian Mixture Model

<sup>3</sup> Morphological Decomposition

<sup>4</sup> Sub-word

<sup>5</sup> Lexicon

پسین است. تعداد نورون‌های خروجی شبکه برابر تعداد کل حالات مدل‌های مخفی مارکوف به کار گرفته شده برای مدل‌سازی شناسه‌ها است که خروجی هر یک از این نورون‌ها احتمال پسین حالت مربوطه به ازای بردار ویژگی ورودی آن است ( $P(s|x)$ ). از آنجا که نیاز به احتمال درست‌نمایی حالات است، احتمال پسین هر حالت بر احتمال پیشین آن تقسیم می‌شود. احتمال پیشین حالات بر مبنای داده‌های آموزشی و به کمک الگوریتم هم‌ترازی اجباری<sup>۱</sup> به دست می‌آید. در رویکرد ترکیبی، شبکه‌های عصبی ژرف مختلفی نظیر TDNN، MLP و LSTM استفاده می‌شود. در این رویکرد برای شروع فرایند آموزش نیاز به وجود یک هم‌ترازی اولیه بین بردارهای ویژگی و حالات HMM است که برای این منظور به طور معمول از یک سامانه مبتنی بر رویکرد HMM-GMM برای ایجاد این هم‌ترازی اولیه بهره گرفته که این وابستگی به نوعی ضعف این رویکرد محسوب می‌شود. از پژوهش‌های اخیر مبتنی بر رویکرد ترکیبی می‌توان از [38] نام برد. این پژوهش که تمرکز اصلی آن بر روی بهبود نرمال‌سازی تصویر و همچنین بررسی روش‌های مختلف استخراج ویژگی است، از شبکه‌های عصبی موجود در نرم‌افزار متن‌باز KALDI برای پیاده‌سازی مدل ترکیبی بهره می‌برد.

در سال ۲۰۱۴ رویکرد انتها به انتها<sup>۲</sup> ارائه شد که مدل مخفی مارکوف به طور کامل کنار گذاشته شده و از یک شبکه عصبی بازگشتی ژرف که به طور معمول از نوع LSTM است به همراه روش CTC<sup>۳</sup> [13] برای مدل‌سازی بهره می‌برد. ورودی شبکه، بردارهای ویژگی است و خروجی آن شناسه‌های زبان به همراه شناسه فاصله. بنابراین تعداد نورون‌های خروجی تعداد شناسه‌های نوشتاری زبان به علاوه یک است. از آنجا که در این رویکرد، شبکه می‌بایست هر دو مدل ویزوال و زبانی را به همراه هم یاد بگیرد، بنابراین باید به اندازه کافی بزرگ باشد و داده آموزشی کافی در اختیار داشته باشیم. به همین خاطر به کارگیری این رویکرد تنها در مواقعی توصیه می‌شود که داده آموزشی زیاد و بستر محاسباتی قوی در اختیار است. از مزایای این رویکرد سادگی ساختار آن و عدم نیاز به مجموعه‌واژگان است که امکان بازشناسی کلمات خارج از مجموعه‌واژگان<sup>۴</sup> را فراهم می‌کند. گفتنی است که عدم نیاز رویکرد انتها به انتها

وجود مجموعه‌واژگان در کنار مزیتی که دارد یک عیب نیز دارد. در دو رویکرد گفته شده قبلی که مبتنی بر وجود مجموعه‌واژگان هستند، فضای جستجو در هنگام رمزگشایی<sup>۵</sup> به کلمات مجموعه‌واژگان محدود شده است و همچنین استفاده از مدل زبانی حین رمزگشایی باعث محدودتر فضای جستجو می‌شود که در نهایت به بهبود دقت بازشناسی منجر می‌شود.

در سال ۲۰۱۷ پژوهش [1] از رویکرد انتها به انتها به کمک روش CTC برای بازشناسی واژگان پیوسته در زبان عربی بهره برد. این پژوهش بر روی دادگان KHATT انجام شده و بازشناسی در سطح شناسه انجام می‌شود. برای این منظور از یک شبکه عصبی بازگشتی LSTM چندجهته<sup>۶</sup> استفاده شده است. پژوهش [17] که در سال ۲۰۱۸ انجام شده است، با بررسی حالات مختلف برای ترکیب چند مدل مبتنی بر شبکه‌های عصبی بازگشتی BLSTM-CTC تلاش می‌کند تا بهترین شیوه ترکیب مدل‌ها را ارائه کند. در این کار ترکیب مدل‌ها در سه سطح پایین، سطح میانی و سطح بالا بررسی می‌شود که نتایج نشان می‌دهد ترکیب سطح بالای مدل‌ها بیشترین بهبود را در نرخ خطا به همراه داشته است. ترکیب سطح بالا به این صورت است که چند مدل BLSTM-CTC به صورت مجزا آموزش داده و خروجی رمزگشایی همه آنها با هم ترکیب می‌شوند.

همان‌طور که در ابتدای این فصل بیان شد، مسأله مدل‌سازی و بازشناسی دست‌نوشته شباهت بسیار زیادی به مسأله مدل‌سازی و بازشناسی گفتار دارد. هر دو مسأله ذاتاً دنباله‌ای هستند و تنها تفاوت آنها فقط در جنس سیگنال ورودی است. به همین علت انتظار می‌رود مدل‌هایی که در مسأله بازشناسی گفتار عملکرد مناسبی داشته‌اند، در مسأله بازشناسی دست‌نوشته نیز عملکرد قابل قبولی داشته باشند. امروزه روش‌های مبتنی بر رویکرد ترکیبی به عنوان مرز دانش در حوزه بازشناسی گفتار شناخته می‌شوند که عمده نقطه ضعف آنها وابستگی به یک هم‌ترازی شناسه‌ای اولیه در شروع روند آموزش است. خوشبختانه مدل‌هایی که در همین اواخر بر مبنای رویکرد ترکیبی ارائه شده‌اند، تلاش کرده‌اند تا وابستگی به این هم‌ترازی اولیه را از بین ببرند که تا حد زیادی نیز موفق بوده‌اند. با توجه به عملکرد مرز دانشی مدل‌های مبتنی بر رویکرد ترکیبی و همچنین سرعت به نسبت بالا آموزش آنها، در این پژوهش این دسته از مدل‌ها را برای مدل‌سازی انتخاب کرده‌ایم.

<sup>5</sup> Decoding

<sup>6</sup> Multi-Directional Long-Short-Term-Memory

<sup>1</sup> Force Alignment

<sup>2</sup> End-to-End

<sup>3</sup> Connectionist Temporal Classification

<sup>4</sup> Out Of Vocabulary



## ۴- اجزای سامانه

این مقاله به کمک روش‌های پیش‌پردازش و نرمال‌سازی تصاویر که در [38] ارائه شده است و همچنین بهره‌گیری از مدل‌سازی‌های به‌روز مبتنی بر رویکرد ترکیبی، سامانه‌ای را ارائه می‌دهد که بهبود چشم‌گیری در میزان دقت بازشناسی کلمات بر روی دادگان عربی KHATT را به همراه داشته است. در ادامه به بررسی جزئیات روش‌های پیش‌پردازش و نرمال‌سازی تصاویر ورودی و روش‌های مدل‌سازی ترکیبی به کار گرفته‌شده خواهیم پرداخت.

### ۴-۱- پیش‌پردازش و نرمال‌سازی تصاویر

پیش‌پردازش و نرمال‌سازی تصاویر ورودی در حوزه بازشناسی دست‌نوشته به‌خصوص دست‌نوشته‌های عربی اهمیت بالایی دارد. پژوهش [38] با تکیه بر نرمال‌سازی تصاویر قبل از آموزش مدل توانسته است بهبود قابل توجهی در میزان دقت به‌دست‌آمده بر روی دادگان عربی KHATT گزارش کند. در ادامه به بررسی جزئیات نرمال‌سازی‌های انجام‌شده در این پژوهش می‌پردازیم.

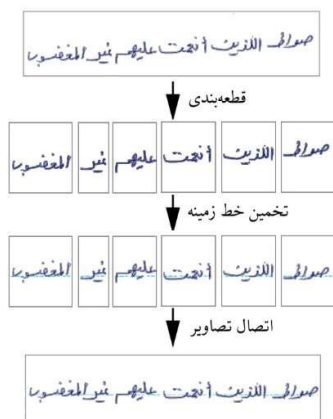
نخستین مرحله در نرمال‌سازی تصویر ورودی، تخمین خط زمینه در محدوده کلمه است. در این جا، این کار توسط روشی که در [37] ارائه شده است، انجام می‌شود. به این صورت که باریک‌ترین نواری که بیشترین بخش از پیکسل‌های یک واژه را در تصویر در بر بگیرد به‌عنوان خط زمینه بالا و پایین انتخاب و در نهایت خط زمینه پایین به‌عنوان خط زمینه لحاظ می‌شود.



(شکل-۳): تخمین خط زمینه بالا و پایین  
(Figure-3): Upper and lower baseline estimation

برای تخمین خط زمینه یک خط دست‌نوشته، ابتدا تصویر به‌صورت عمودی به بخش‌هایی که هر کدام شامل یک کلمه هستند تقسیم می‌شود؛ سپس برای هر بخش که شامل یک کلمه است به روش گفته‌شده در بالا، خط زمینه تخمین زده می‌شود. بعد از تخمین خط زمینه، تصویر را حول این خط به میزانی می‌چرخانیم که خط زمینه به‌طور کامل افقی شود؛ سپس همه بخش‌ها را به هم متصل می‌کنیم تا تصویر

نرمال‌شده کل خط به‌دست آید. این کار کجی تصاویر خطوط را اصلاح می‌کند.



(شکل-۴): مراحل تخمین و تصحیح خط زمینه  
(Figure-4): Steps of baseline estimation and correction

یکی دیگر از ویژگی‌های تصاویر دست‌نوشته‌های عربی این است که حروف به‌طور معمول به‌طور کامل در راستای عمودی قرار ندارند و با توجه به سبک دست‌خط نویسنده قدری به چپ و یا راست متمایل هستند. مرحله بعدی در نرمال‌سازی تصویر، از بین بردن این اریبی کلمات است. یکی از روش‌ها برای از بین بردن این اریبی استفاده از تبدیل Shear است. برای پیدا کردن زاویه مناسب در این تبدیل از تبدیل Hough [10] مطابق آنچه در [37] آمده استفاده شده است.



(شکل-۵): اریبی‌های متفاوت در دست‌نوشته فارسی  
(Figure-5): Slants in Persian Handwritten

در این پژوهش از پیکسل‌های خام تصویر به‌عنوان بردار ویژگی استفاده شده است. به این صورت که یک پنجره به عرض یک پیکسل و همچنین با طول گام یک پیکسل بر روی تصویر لغزانیده می‌شود و یک بردار ویژگی را به‌اندازه عرض تصویر که مقادیر خام پیکسل‌ها هستند، تولید می‌کند. به این دلیل که ضروری است، تعداد ویژگی‌های تمام بردارهای ویژگی با یکدیگر برابر باشد؛ به این منظور بعد از مراحل نرمال‌سازی، تمام تصاویر را با حفظ نسبت طول به عرض آنها، به تصاویری با عرض پنجاه پیکسل تغییر اندازه می‌دهیم.

در روند آموزش مدل های مبتنی بر رویکرد ترکیبی (IIMM-DNN) که هدف این پژوهش است، می توان از توابع هدف مختلفی استفاده کرد. این توابع هدف را می توان به دو دسته توابع هدف در سطح قاب و توابع هدف در سطح دنباله تقسیم کرد. در مدل هایی که توابع هدف آنها در سطح قاب تعریف شده است، به روزرسانی پارامترهای مدل بعد از مشاهده هر بردار ویژگی از ورودی انجام می شود؛ در حالی که در مدل هایی که توابع هدف آنها دنباله ای است، ابتدا باید تمام بردارهای ویژگی داده آموزشی مشاهده و سپس پارامترهای مدل به روزرسانی شوند. از آن جا که بازنمایی دست نوشته یک مسئله ذاتاً دنباله ای است، انتظار می رود توابع هدف دنباله ای عملکرد بهتری داشته باشند. از جمله مهم ترین توابع دنباله ای معرفی شده در حوزه بازنمایی گفتار که می توان در بازنمایی دست نوشته نیز از آنها بهره برد،  $MMI^1$  [4]،  $sMBR^2$  [12]،  $MWE^3$  و  $MPE^4$  [28] هستند که شباهت های بسیار زیادی به یکدیگر دارند. پژوهش های اخیر در حوزه بازنمایی گفتار نشان داده [40,39,35] که با استفاده از این توابع بهبود در خور توجهی قابل حصول است که انتظار می رود به کارگیری آنها در بازنمایی دست نوشته نیز مؤثر واقع شود.

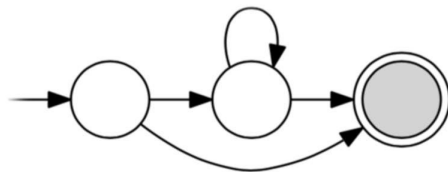
پایه توابع هدف تمایزی بر مبنای بیشینه سازی اطلاعات متقابل بین  $M_{w(u)}$  به عنوان گراف آموزشی و  $x^{(u)}$  به عنوان دنباله بردارهای ویژگی ورودی نام به ازای تمام ورودی ها است که به صورت زیر تعریف می شود:

$$\mathcal{F}_{MMI} = \sum_{u=1}^U I(x^{(u)}, M_{w(u)}) = \sum_{u=1}^U \log \frac{p_{\lambda}(x^{(u)} | M_{w(u)})}{p_{\lambda}(x^{(u)})} \quad (2)$$

که در آن  $I$  تابع اطلاعات مشترک و  $U$  تعداد کل نمونه های آموزشی است؛ بنابراین این تابع هدف که تابع هدف MMI نامیده می شود، علاوه بر این که در سطح دنباله است، برخلاف تابع هدف بیشینه-شباهت یک تابع هدف تمایزی نیز است. یک رویکرد ترکیبی که از این تابع هدف بهره برده است LF-MMI<sup>5</sup> [27] نامیده می شود و در حال حاضر مرز دانش در

بازشناسی گفتار به حساب می آید. روش LF-MMI در [14] به نحوی تغییر داده شده است که نیاز به هم ترازی قبلی ندارد. این روش جدید LF-MMI تخت آغاز<sup>6</sup> نامیده شده است.

روش LF-MMI در تئوری تفاوتی با روش MMI ندارد. محاسبه بر روی گراف مخرج در روش MMI مسأله چالش برانگیزی است. به دلیل این که گراف مخرج، گرافی است متشکل از تمام حالات ممکن، بنابراین محاسبات بر روی آن بسیار پر هزینه است؛ از این رو به طور معمول محاسبه گراف مخرج به وسیله  $n$ -تا بهترین<sup>7</sup> و یا لاتیس تقریب زده می شود. به این معنی که به جای محاسبه بر روی تمام حالات ممکن، محاسبات تنها بر روی محتمل ترین حالات انجام می شود. در روش MMI همچنین برای مقداردهی اولیه شبکه عصبی از پیش آموزش شبکه به وسیله شبکه های باور عمیق<sup>8</sup> استفاده می شود [7].



(شکل-۶): توپولوژی شناسه در روش LF-MMI  
(Figure-6): HMM topology for characters in LF-MMI

تفاوت روش LF-MMI با روش آموزش تمایزی MMI در این است که گراف مخرج به جای  $n$ -تا بهترین و یا لاتیس به صورت دقیق محاسبه و همچنین آموزش شبکه نیز با مقداردهی اولیه تصادفی شروع می شود. برای کاهش هزینه محاسبات گراف مخرج سه روش پیشنهاد شده است:

۱- در گراف مخرج به جای استفاده از مدل زبانی در سطح واژه، از مدل زبانی در سطح شناسه استفاده می شود. این کار باعث می شود اندازه گراف دو تا سه مرتبه کوچک تر شود.

۲- جهت انجام الگوریتم های پیش رو-پس رو<sup>9</sup> برای گراف مخرج از پردازنده های گرافیکی استفاده می شود.

۳- نرخ قاب در خروجی شبکه به یک سوم کاهش داده می شود. این کار با استفاده از شبکه های TDNN به سادگی قابل انجام است. این کار باعث می شود تا هم

<sup>6</sup> Flat start

<sup>7</sup> N-best

<sup>8</sup> Deep Belief Network (DBN)

<sup>9</sup> Forward-Backward algorithms

<sup>1</sup> Maximum Mutual Information

<sup>2</sup> State-level Minimum Bayes Risk

<sup>3</sup> Minimum Word Error

<sup>4</sup> Minimum Phone Error

<sup>5</sup> Lattice-Free Maximum Mutual Information



گره نزیم، دقت به‌دست آمده تغییر چندانی نمی‌کند؛ زیرا شبکه یاد می‌گیرد تا خروجی حالت‌های بدون استفاده را همیشه صفر تولید کند و این باعث می‌شود تا این حالت‌ها در هیچ کدام از محاسبات گراف صورت و مخرج اثری نداشته باشند.

این مقاله به‌کمک روش‌های پیش‌پردازش و نرمال‌سازی تصاویر که در [38] ارائه شده است و همچنین بهره‌گیری از رویکردهای ترکیبی LF-MMI و LF-MMI تحت‌آغاز جهت مدل‌سازی، سامانه‌ای را ارائه می‌دهد که بهبود چشم‌گیری در میزان دقت بازشناسی واژگان بر روی دادگان عربی KHATT را به‌همراه داشته است. شبکه عصبی استفاده‌شده در کلیه آزمایش‌های انجام‌شده در این مقاله شامل هفت لایه پیچشی<sup>۲</sup> در ابتدا و پنج لایه تأخیر زمانی<sup>۳</sup> در انتهاست.

## ۵- نتایج و آزمایش‌ها

آزمایش‌ها بر روی دادگان عربی KHATT انجام شده است. دادگان KHATT شامل تصاویر خطوط دست‌نوشته است که از فرم‌های تکمیل شده توسط هزار نویسنده استخراج شده است. این دادگان شامل دو بخش مجزا است. بخش نخست شامل جملاتی ثابتی است که توسط تمام نویسنده‌ها نوشته شده و شامل تمام شناسه‌ها و حالات مختلف آنها است. بخش دوم، شامل جملاتی یکتاست که از ۴۶ منبع مختلف جمع‌آوری شده و برای هر نویسنده متفاوت است. این دادگان در مجموع شامل ۱۳۳۵۷ تصویر خطوط دست‌نوشته است که از بین آنها ۹۴۶۲ خط برای مجموعه آموزش، ۱۸۹۹ خط برای مجموعه توسعه و ۱۹۹۶ خط برای مجموعه آزمون در نظر گرفته شده است. در این پژوهش برای تقسیم‌بندی تصاویر به مجموعه‌های آموزش، توسعه و آزمون مشابه [16] و [38] عمل شده است تا نتایج قابل مقایسه باشد. در انجام آزمایش‌ها، از یک مدل زبانی 3-gram استفاده شده که از متون داده آموزشی استخراج شده است. نرخ کلمات خارج از مجموعه‌واژگان برای این دادگان برابر ۱۱/۴ درصد است. عملکرد سامانه بر مبنای درصد خطای بازشناسی واژه ارزیابی شده است که مطابق رابطه زیر محاسبه می‌شود:

$$(3) \quad 100 * \frac{\text{تعداد جایگزینی} + \text{تعداد حذف} + \text{تعداد درج}}{\text{تعداد کل کلمات}}$$

پیچیدگی محاسبات در خود شبکه کاهش و همچنین تعداد دفعاتی که پارامترهای شبکه را به‌روزرسانی می‌کنیم، کاهش یابد. کاهش نرخ قاب در خروجی باعث می‌شود تا به‌ازای هر سه قاب از ورودی، یک بردار را از احتمال‌های پسین در خروجی داشته باشیم. به همین علت باید توپولوژی مدل مارکوف شناسه‌ها را به‌نحوی تغییر دهیم که با یک گذر بتوان از هر شناسه عبور کرد؛ در غیر این صورت کوتاه‌ترین طول برای شناسه سه برابر حالت عادی می‌شود که مسلماً تأثیر منفی بر دقت بازشناسی شناسه‌ها خواهد داشت. توپولوژی HMM استفاده‌شده در روش LF-MMI در شکل (۶) مشاهده می‌شود.

ضعف عمده مدل LF-MMI مشابه سایر رویکردهای ترکیبی این است که آموزش آن نیاز به یک هم‌ترازی اولیه دارد. هدف اصلی از ارائه مدل LF-MMI تحت‌آغاز مرتفع‌کردن این مشکل است. به عبارتی این روش در تلاش است که تمام فرایند آموزش در یک مرحله انجام شود و نیازی به وجود هم‌ترازی اولیه نباشد. روش LF-MMI در دو مرحله از یک هم‌ترازی اولیه استفاده می‌کند. نخستین مرحله برای ساخت گراف صورت است. گراف صورت در این روش بر مبنای لاتیس ساخته‌شده به‌وسیله یک مدل آموزش‌دیده‌شده، طول شناسه‌ها را تخمین می‌زند و گرافی تولید می‌کند که در آن هیچ طوقه‌ای<sup>۱</sup> وجود ندارد. وابستگی بعدی در مرحله آموزش درخت تصمیم بر مبنای هم‌ترازی اولیه است که به‌وسیله آن حالت‌های HMM به هم گره زده می‌شوند. در روش LF-MMI تحت‌آغاز برای حذف این وابستگی در مرحله نخست، از یک گراف ترکیبی که برای آموزش مدل‌های HMM-GMM استفاده می‌شود، به‌عنوان گراف صورت استفاده می‌کنیم. این گراف هیچ محدودیتی برای داشتن طوقه ندارد؛ از این‌رو نیازی به هیچ هم‌ترازی اولیه‌ای نیز ندارد. برای حذف وابستگی درخت تصمیم مورد استفاده به هم‌ترازی اولیه در روش LF-MMI تحت‌آغاز، از تمام ترکیب‌های دوشناسه‌ای استفاده شده است؛ درواقع هیچ کدام از حالت‌های HMM به هم گره زده نمی‌شوند. در آزمایش‌های [27] نشان داده شده است که استفاده از ترکیب‌های دوشناسه‌ای عملکردی معادل ترکیب‌های سه‌شناسه‌ای دارند؛ همچنین در آزمایش‌های اخیر [14] نشان داده شده است که اگر هیچ کدام از حالت‌ها را به هم

<sup>1</sup> Loop

<sup>2</sup> Convolutional

<sup>3</sup> Time Delay

(جدول ۳-): درصد خطای بازشناسی کلمه برای مجموعه توسعه و آزمون تعریف شده بر روی دادگان khatt به ازای روش های مختلف

استخراج ویژگی و مدل سازی در نرم افزار KALDI

(Table-3): Word error rate of recognition on the test set and development set of KHATT data set for various feature extraction and modeling techniques in KALDI

شماره	ویژگی	واحد مدل سازی	مدل سازی	خطای بازشناسی کلمه (%)	
				مجموعه توسعه	مجموعه آزمون
1	Raw pixels + delta	Mono-Character	HMM-GMM	103.87	103.41
2	Raw pixels + delta	Tri-Character	HMM-GMM	78.05	76.46
3	Raw pixels + LDA + MLLT	Tri-Character	HMM-GMM	56.08	54.14
4	Raw pixels + LDA + MLLT + WAT	Tri-Character	HMM-GMM	55.46	53.21
5	Raw pixels	Bi-Character	HMM-DNN (LF-MMI)	29.31	27.96

دارند، در این مدل، ابتدا به کمک مدل آموزش داده شده در سطر چهارم، یک هم ترازوی اولیه برای بردارهای ویژگی داده های آموزشی به دست می آوریم و از آن در روند آموزش شبکه عصبی رویکرد ترکیبی LF-MMI استفاده می کنیم. نتایج نشان می دهند که استفاده از این مدل باعث بهبود مطلق ۲۵/۲۵ درصدی در خطای بازشناسی کلمات شده است.

#### ۱-۵- بررسی تأثیر هم ترازوی شناسه های اولیه

همان طور که در بخش ۳ بیان شد، مدل LF-MMI تخت آغاز گسترشی از روش LF-MMI است که نیاز به هم ترازوی شناسه های اولیه ندارد. در سطر نخست جدول (۴) نتایج آزمایش با استفاده از مدل LF-MMI تخت آغاز قابل مشاهده است. همچنین با توجه به این که مدل LF-MMI نیازمند یک هم ترازوی اولیه است، این مدل را دوباره و با استفاده از هم ترازوی های اولیه متفاوت مورد آزمایش قرار دادیم. در سطر دوم جدول (۴)، از هم ترازوی به دست آمده از مدل سطر چهارم جدول (۳) استفاده شده است. در سطر سوم، هم ترازوی اولیه به دست آمده به وسیله مدل سطر پنجم جدول (۳) که خود که یک مدل ترکیبی است، استفاده شده است. در انتها و در سطر چهارم، هم ترازوی به دست آمده به وسیله مدل LF-MMI تخت آغاز، به عنوان هم ترازوی شناسه های اولیه مورد استفاده قرار گرفته است. با توجه به نتایج به دست آمده، می توان مشاهده کرد که هرچه هم ترازوی اولیه دقیق تر باشد، دقت روش LF-MMI نیز بیش تر است؛ همچنین دقتی که روش LF-MMI تخت آغاز به همراه داشته با توجه به بی نیاز بودن از هم ترازوی اولیه، دقت درخور توجهی است.

در جدول (۳) نتایج بازشناسی بر روی مجموعه توسعه و آزمون دادگان KHATT به ازای ویژگی ها و مدل های مختلف آمده است. سطر نخست، نتایج را برای ساده ترین حالت، یعنی استفاده از پیکسل های خام و مشتق نخست آنها به عنوان ویژگی و همچنین بهره بردن از مدل مبتنی بر تک شناسه نشان می دهد. در سطر دوم تفاوت نسبت به سطر نخست تنها در روش مدل سازی اتفاق افتاده است که مدل سازی تک شناسه ای به سه شناسه ای تغییر یافته است و همان طور که مشاهده می شود، بهبود قابل توجهی به همراه داشته است. در سطر سوم، بر روی هر بردار ویژگی که در واقع یک ستون از پیکسل های تصویر اصلی است، به همراه بردارهای ویژگی قبل و بعد، دو تبدیل خطی <sup>۱</sup>LDA و <sup>۲</sup>MLLT به ترتیب اعمال شده و در نهایت یک بردار چهل مؤلفه ای به دست می آید که مدل سازی سه شناسه ای بر مبنای آن صورت می گیرد. همان طور که در جدول نیز مشاهده می شود، استفاده از تبدیل های LDA و MLLT به جای مشتق نخست، باعث بهبود مطلق خطای بازشناسی به میزان ۲۲/۳۲ درصد شده است. با توجه به این بهبود قابل توجه، از این مرحله به بعد در تمام مراحل از این دو تبدیل استفاده شده است. در سطر چهارم، آموزش تطبیقی به نویسندگان<sup>۳</sup> صورت گرفته است. همچنین در مرحله بازشناسی، ویژگی های هر گوینده به کمک روش IMLLR تطبیق داده شده که در نهایت منجر به کاهش ۰/۹۳ درصدی خطای بازشناسی شده است. در سطر پنجم تنها از پیکسل های خام تصویر به عنوان بردار ویژگی استفاده شده است. به این علت که مدل های ترکیبی نیاز به یک هم ترازوی شناسه ای اولیه

<sup>۱</sup> Linear Discriminant Analysis

<sup>۲</sup> Maximum Likelihood Linear Transform

<sup>۳</sup> Writer Adaptive Training (WAT)

(جدول-۴): نتایج به‌دست‌آمده از دو هم‌ترازی اولیه متفاوت

(Table-4): Results of two different initial alignments

مدل‌سازی	هم‌ترازی اولیه	خطای بازشناسی واژه (%)	
		توسعه	آزمون
HMM-DNN (Flatstart LF-MMI)	-	28.87	26.83
HMM-DNN (LF-MMI)	HMM-GMM	29.31	27.96
HMM-DNN (LF-MMI)	HMM-DNN (LF-MMI)	27.66	26.26
HMM-DNN (LF-MMI)	HMM-DNN (Flatstart LF-MMI)	26.21	25.32

مدل‌های زبانی مختلف به‌وسیله مدل LF-MMI با هم‌ترازی اولیه به‌دست‌آمده از مدل تخت‌آغاز مورد آزمایش قرار گرفتند. این مدل‌های شامل 2-gram و 3-gram بودند که هر کدام به‌وسیله روش‌های هموارسازی kneser-good-turing، kneser-ney و بیشینه بی‌نظمی<sup>۸</sup> هموار شدند. جدول (۶) نتایج به‌دست‌آمده از این آزمایش‌ها را نشان می‌دهد.

(جدول-۶): نتایج به‌دست‌آمده از انواع مدل زبانی و روش‌های

مختلف هموارسازی

(Table-6): Results of various language model types and smoothing techniques

مدل	هموارسازی	خطای بازشناسی کلمه (%)	
		توسعه	آزمون
2-gram	Good-Turing	25.20	24.00
2-gram	Kneser-Ney	24.57	23.19
2-gram	MaxEnt	24.96	23.68
3-gram	Good-Turing	25.34	24.14
3-gram	Kneser-Ney	24.88	23.52
3-gram	MaxEnt	24.60	23.18

همان‌طور که در جدول مشاهده می‌شود، استفاده از مدل زبانی 2-gram باعث کاهش نرخ خطا شده است. این کاهش به ازای استفاده از روش هموارسازی Kneser-Ney به بیشترین حد خود رسیده و همچنین استفاده از روش‌های هموارسازی Kneser-Ney و بیشینه بی‌نظمی در مدل 3-gram نیز باعث کاهش نرخ خطا شده است. در مدل 3-gram، روش هموارسازی بیشینه بی‌نظمی بیشترین بهبود را داشته است.

#### ۴-۵- مقایسه نتایج با سایر پژوهش‌ها

در جدول (۷) نتایج به‌دست‌آمده در این پژوهش با سایر پژوهش‌ها مقایسه شده است. تا جایی که مطلع هستیم بهترین دقت بر روی دادگان خط با استفاده از مدل زبانی استخراج‌شده از متن داده آموزشی توسط [38] ارائه شده است. هر دو پژوهش [16] و [38] از مدل زبانی 3-gram استخراج‌شده از متون داده آموزشی استفاده می‌کنند. به همین علت ما نیز بهترین دقتی را که به‌وسیله یک مدل زبانی 3-gram حاصل شده با این پژوهش‌ها مقایسه می‌کنیم. همان‌طور که در جدول (۷) قابل مشاهده است، دقت بازشناسی واژه به‌دست‌آمده در این پژوهش بهترین دقت گزارش‌شده را به میزان ۷/۳۲ درصد مطلق بهبود بخشیده است.

<sup>8</sup> Maximum Entropy

#### ۲-۵- بررسی اثر افزایش داده<sup>۱</sup>

جهت بررسی اثر افزایش داده بر نرخ خطای بازشناسی کلمه، با انجام تبدیلات مختلفی از جمله افزودن نویز گاوسی، ایجاد ارببی<sup>۲</sup>، چرخش<sup>۳</sup>، حذف تصادفی<sup>۴</sup> پیکسل‌ها، مات کردن<sup>۵</sup> و تغییر مقیاس<sup>۶</sup> بر روی تصاویر آموزشی، تلاش شد تا داده‌های آموزشی را به‌صورت مصنوعی افزایش دهیم. تبدیل‌های انجام‌شده باعث می‌شود تا مدل در هنگام آموزش تصاویر متنوعی ببیند که این باعث کاهش بیش‌برازش<sup>۷</sup> و افزایش قابلیت تعمیم‌پذیری مدل می‌شود.

با آموزش مدل ترکیبی LF-MMI با استفاده از هم‌ترازی به‌دست‌آمده توسط مدل LF-MMI تخت‌آغاز به‌عنوان هم‌ترازی اولیه، بر روی داده‌های آموزشی افزایش داده‌شده، نتایجی که در جدول (۵) مشاهده می‌شود به‌دست آمد. نتایج نشان می‌دهند که نرخ خطای بازشناسی کلمه بر روی داده توسعه به میزان ۰/۸۷ درصد و بر روی داده آزمون به میزان ۱/۱۸ درصد بهبود مطلق داشته است.

(جدول-۵): نتایج به‌دست‌آمده از آموزش مدل بر روی داده

آموزشی افزایش یافته

(Table-5): Results of model training using augmented training data

مدل	نرخ خطای واژه (%)	
	توسعه	آزمون
LF-MMI (without augmentation)	26.21	25.32
LF-MMI (with augmentation)	25.34	24.14

#### ۳-۵- بررسی اثر مدل زبانی

از دیگر آزمایش‌های انجام‌شده در این پژوهش، بررسی اثر انواع مختلف مدل زبانی بر روی نتایج به‌دست‌آمده است.

<sup>1</sup> Data Augmentation

<sup>2</sup> Skewing

<sup>3</sup> Rotating

<sup>4</sup> Dropout

<sup>5</sup> Blurring

<sup>6</sup> Scaling

<sup>7</sup> Over-fitting

## 7- References

## ۷- مراجع

- [1] R. Ahmad, S. Naz, M. Z. Afzal, S. F. Rashid, M. Liwicki, A. Dengel, "KHATT: A Deep Learning Benchmark on Arabic Script.", In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, pp. 10-14, 2017.
- [2] R. Al-Hajj, C. Mokbel, C. Mokbel, L. Likforman-Sulem and L. Likforman-Sulem, "Combination of HMM-Based Classifiers for the Recognition of Arabic Handwritten Words", Ninth International Conference on Document Analysis and Recognition (ICDAR 2007) Vol 2, 2007.
- [3] J. AlKhatcb, J. Ren, J. Jiang and H. Al-Muhtaseb, "Offline handwritten Arabic cursive text recognition using Hidden Markov Models and re-ranking", *Pattern Recognition Letters*, vol. 32, no. 8, pp. 1081-1088, 2011.
- [4] L. Bahl, P. Brown, P. de Souza and R. Mercer, "Maximum mutual information estimation of hidden Markov model parameters for speech recognition", in *International Conference on Acoustics, Speech, and Signal Processing*, vol. 11, pp. 49-52, 1986.
- [5] A. M. Bidgoli, M. Sarhadi, "IAUT/PHCN: Islamic Azad University of Tehran/Persian handwritten city names, a very large database of handwritten Persian word.", *11th International Conference on Frontiers in Handwriting Recognition*, pp. 192-197, 2008.
- [6] A. Broumandnia, J. Shanbehzadeh and M. Reza khah Varnoosfaderani, "Persian/arabic handwritten word recognition using M-band packet wavelet transform", *Image and Vision Computing*, vol. 26, no. 6, pp. 829-842, 2008.
- [7] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large- vocabulary speech recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 1, pp. 30-42, 2012.
- [8] M. Dchghan, K. Facz, M. Ahmadi and M. Shridhar, "Unconstrained Farsi handwritten word recognition using fuzzy vector quantization and hidden Markov models", *Pattern Recognition Letters*, vol. 22, no. 2, pp. 209-214, 2001.
- [9] M. Dchghan, K. Facz, M. Ahmadi, M. Shridhar, "Handwritten Farsi (Arabic) word recognition: a holistic approach using discrete HMM," *Pattern Recognition*, vol. 34, no. 5, pp. 1057-1065, 2001.
- [10] R. Duda and P. Hart, "Use of the Hough transformation to detect lines and curves in pictures", *Communications of the ACM*, vol. 15, no. 1, pp. 11-15, 1972.
- [11] A. Elbaati, H. Boubaker, M. Kherallah, A. Ennaji, H. Abed and A. Alimi, "Arabic Handwriting Recognition Using Restored Stroke

(جدول-۷): مقایسه نتایج به دست آمده با سایر نتایج گزارش شده

بر روی دادگان KHATT

(Table-7): Comparison of our results with other reported results on KHATT dataset

پژوهش	نرخ خطای واژه (%)	
	توسعه	آزمون
Hamdani et al. [16]	33.60	34.10
QCRI [38]	29.40	30.50
Jemni et al. [17]	31.60	29.44
پژوهش جاری	24.60	23.18

## ۶- جمع بندی

در این مقاله بعد از مروری اجمالی بر روی پژوهش‌های صورت گرفته و منابع موجود در حوزه بازشناسی دست‌نوشته عربی و فارسی، به بررسی شباهت‌های دو مسأله بازشناسی گفتار و بازشناسی دست‌نوشته پرداختیم؛ سپس با مروری کامل بر روی رویکردهای مرسوم برای مدل‌سازی و پژوهش دست‌نوشته، با توجه به نتایج قابل قبول رویکردهای مبتنی بر HMM-DNN در حوزه بازشناسی گفتار و همچنین بی‌نیازی این رویکرد از انجام استخراج ویژگی صریح و مهندسی‌شده از تصاویر ورودی به دلیل وجود شبکه عصبی، این رویکرد را برای آزمایش بر روی دادگان KHATT انتخاب کردیم. با به کارگیری مدل‌های LF-MMI و LF-MMI تخت‌تخت‌آغاز که از جدیدترین مدل‌های مبتنی بر HMM DNN هستند، و همچنین بهره‌بردن از روش افزایش داده موفق به بهبود نرخ خطای کلمه بر روی این دادگان شدیم. بهترین نرخ خطای واژه بر روی مجموعه توسعه برابر با ۲۴/۶۰ و بر روی مجموعه آزمون برابر با ۲۳/۱۸ به دست آمد که بهبود درخور توجهی نسبت به نتایج گزارش شده قبلی بر روی دادگان KHATT است. تمرکز اصلی این پژوهش بر روی به کارگیری مدل‌های مرسوم در حوزه بازشناسی گفتار برای حل مسأله بازشناسی دست‌نوشته عربی بوده است.

## سپاس‌گزاری

نویسندگان مقاله برخود لازم می‌دانند از همکاری و حمایت آقای دبلیو پژوهش‌گر ارشد مرکز پردازش زبان و گفتار (CLSP) دانشگاه جانز هاپکینز و دکتر حسین هادیان فارغ‌التحصیل مقطع دکترای هوش مصنوعی دانشگاه صنعتی شریف به خاطر راهنمایی‌های بی‌دریغ‌شان و فهم عمیق عمادالدین پژوهش‌گر مؤسسه پژوهشی QCRI کشور قطر به خاطر به اشتراک گذاشتن برخی منابع و کدها صمیمانه سپاس‌گزاری کند.

- of Farsi handwritten city names." In International Conference on Frontiers in Handwriting Recognition. 2008.
- [23] S. Mozaffari, K. Faez, F. Faradji, M. Ziaratban, S. M. Golzan, "A comprehensive isolated Farsi/Arabic character database for handwritten OCR research," In Tenth International Workshop on Frontiers in Handwriting Recognition, Suvisoft, 2006.
- [24] S. Mozaffari, K. Faez and M. Ziaratban, "Structural decomposition and statistical description of Farsi/Arabic handwritten numeric characters", Eighth International Conference on Document Analysis and Recognition (ICDAR'05), 2005.
- [25] M. Pechwitz, S. S. Maddouri, V. Märgner, N. Ellouze, H. Amiri, "IFN/ENIT-database of handwritten Arabic words," In *Proc. of CIFED*, vol. 2, pp. 127-136. 2002.
- [26] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, J. Silovsky, G. Stemmer, K. Vesely, "The Kaldi speech recognition toolkit," In *IEEE 2011 workshop on automatic speech recognition and understanding, IEEE Signal Processing Society*, 2011.
- [27] D. Povey, V. Peddinti, D. Galvez, P. Ghahmani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, "Purely sequence-trained neural networks for asr based on lattice-free mmi," in Interspeech, 2016.
- [28] D. Povey and P. C. Woodland, "Minimum phone error and i-smoothing for improved discriminative training," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP, vol. 1. IEEE*, pp. 1-105, 2002.
- [29] D. Rybach, S. Hahn, P. Lehnen, D. Nolden, M. Sundermeyer, Z. Tüske, S. Wiesler, R. Schlüter, H. Ney, "Rasr-the rwth aachen university open source speech recognition toolkit," In *Proc. IEEE Automatic Speech Recognition and Understanding Workshop*. 2011.
- [30] R. Sabzi, Z. Fotoohinya, A. Khalili, S. Golzari, Z. Salkhorde, S. Behraves, S. Akbarpour, "Recognizing Persian handwritten words using deep convolutional networks," in *Artificial Intelligence and Signal Processing Conference (AISP)*, pp. 85-90, 2017.
- [31] J. Sadri, C. Y. Sucn, and T. D. Bui, "Application of support vector machines for recognition of handwritten Arabic/Persian digits," In *Proceedings of Second Iranian Conference on Machine Vision and Image Processing*, vol. 1, pp. 300-307. 2003.
- Chronology", 2009 10th International Conference on Document Analysis and Recognition, 2009.
- [12] V. Goel and W. Byrne, "Minimum Bayes-risk automatic speech recognition", *Computer Speech & Language*, vol. 14, no. 2, pp. 115-135, 2000.
- [13] A. Graves, S. Fernández, F. Gomez and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," In *Proceedings of the 23rd international conference on Machine learning*, pp. 369-376, 2006.
- [14] H. Hadian, H. Sameti, D. Povey and S. Khudanpur, "Flat-Start Single-Stage Discriminatively Trained HMM-Based Models for ASR", *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 11, pp. 1949-1961, 2018.
- [15] P. Haghighi, N. Nobile, C. He and C. Suen, "A New Large-Scale Multi-purpose Handwritten Farsi Database", *Lecture Notes in Computer Science*, pp. 278-286, 2009.
- [16] M. Hamdani, A. Mousa and H. Ney, "Open Vocabulary Arabic Handwriting Recognition Using Morphological Decomposition", 2013 12th International Conference on Document Analysis and Recognition, 2013.
- [17] S. K. Jemni, Y. Kessentini, S. Kanoun, J. Ogier, "Offline Arabic Handwriting Recognition Using BLSTMs Combination," In *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*, pp. 31-36, 2018.
- [18] S. Khorashadizadeh, A. Latif, "Arabic/Farsi Handwritten Digit Recognition usin Histogram of Oriented Gradient and Chain Code Histogram", *International Arab Journal of Information Technology (IAJIT)*, vol. 13, no. 4, 2016.
- [19] H. Khosravi, E. Kabir, "Introducing a vcrly large dataset of handwritten Farsi digits and a study on their varieties.", *Pattern recognition letters*, vol. 28, no. 10, pp. 1133-1141. 2007.
- [20] D. Lee, S. Ismael, S. Grimes, D. Doermann, S. Strassel, Z. Song, "MADCAT Phase 1 Training Set", LDC2012T15. DVD. Philadelphia: Linguistic Data Consortium, 2012.
- [21] S. A. Mahmoud, I. Ahmad, M. Alshayeb, W. G. Al-Khatib, M. T. Parvez, G. A. Fink, V. Margner, H. E. Abed, "KHATT: Arabic offline handwritten text database," In *2012 International Conference on Frontiers in Handwriting Recognition (ICFHR 2012)*, pp. 449-454, 2012.
- [22] S. Mozaffari, H. E. Abed, V. Märgner, K. Faez, A. Amirshahi, "IfN/Farsi-Database: a database



گرایان تصویر "مجله علمی و پژوهشی پردازش علائم و داده‌ها. ۱۳۹۴؛ ۱۲ (۳): ۱۵-۲۹

- [43] E. Bayesteh Tashk, A. Ahmadifard and H. khosravi, "A two step method for offline handwritten Farsi word recognition using adaptive division of gradient image", *JSDP*, Vol.12 (3), pp.15-29, 2015.



**باقر باباعلی** مدرک کارشناسی خود را در سال ۱۳۷۹ در رشته مهندسی رایانه (گرایش سخت‌افزار) از دانشگاه شیراز و همچنین مدرک کارشناسی ارشد و دکترای خود را به ترتیب در سال‌های ۱۳۸۲ و

۱۳۸۹ در درشته مهندسی رایانه (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 رایانه دانشگاه تهران است.

نشانی رایانامه ایشان عبارت است از:

**babaali@ut.ac.ir**



**بابک رکابدار** مدرک کارشناسی خود را در سال ۱۳۹۵ در رشته علوم رایانه از دانشگاه امیرکبیر دریافت کرد. همچنین مهر ماه سال ۱۳۹۵ به‌منظور دریافت مدرک کارشناسی ارشد (گرایش هوش

مصنوعی) وارد دانشگاه تهران شد. زمینه‌های پژوهشی مورد علاقه ایشان پردازش زبان طبیعی، یادگیری ماشین، بازشناسی دست‌نوشته و یادگیری ژرف است.

نشانی رایانامه ایشان عبارت است از:

**Babak.rekabdar@ut.ac.ir**

- [32] J. Sadri, M. R. Yeganehzad, J. Saghi, "A novel comprehensive database for offline Persian handwriting recognition.", *Pattern Recognition*, vol. 60, pp. 378-393, 2016.
- [33] R. Safabakhsh, A. Ghanbarian and G. Ghiasi, "HaFT: A handwritten Farsi text database", 2013 8th Iranian Conference on Machine Vision and Image Processing (MVIP), 2013.
- [34] H. Sajedi, "Handwriting recognition of digits, signs, and numerical strings in Persian", *Computers & Electrical Engineering*, vol. 49, pp. 52-65, 2016.
- [35] H. Sak, O. Vinyals, G. Heigold, A. Senior, E. McDermott, R. Monga, and M. Mao, "Sequence discriminative distributed training of long short-term memory recurrent neural networks," in *Interspeech*, 2014.
- [36] H. Soltanzadeh, M. Rahmati, "Recognition of Persian handwritten digits using image profiles of multiple orientations," *Pattern Recognition Letters*, vol. 25, no. 14, pp. 1569-1576, 2004.
- [37] F. Stahlberg and S. Vogel, "Detecting dense foreground stripes in Arabic handwriting for accurate baseline positioning", 2015 13th International Conference on Document Analysis and Recognition (ICDAR), 2015.
- [38] F. Stahlberg and S. Vogel, "The QCRI Recognition System for Handwritten Arabic", *Image Analysis and Processing*, pp. 276-286, 2015.
- [39] P. Voigtlaender, P. Doetsch, S. Wiesler, R. Schlüter, and H. Ney, "Sequence-discriminative training of re-current neural networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 2100-2104, 2015.
- [40] K. Vesely, A. Ghoshal, L. Burget, and D. Povey, "Sequence-discriminative training of deep neural networks," in *INTERSPEECH*, pp. 2345-2349, 2013.
- [41] S. Young, G. Evermann, D. Kershaw, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, P. Woodland, *The HTK Book (for version 3.4)*. Cambridge Univ. Eng. Dept., 2009.
- [42] M. Ziaratban, K. Faez and F. Bagheri, "FHT: An Unconstraint Farsi Handwritten Text Database", 2009 10th International Conference on Document Analysis and Recognition, pp. 281-285, 2009.

[43] بایسته، تاشک الهام، احمدی فرد، علیرضا و خسروی، حسین "یک روش دو مرحله‌ای برای بازشناسی کلمات دست‌نوشته فارسی به کمک بلوک‌بندی تطبیقی