

ParsAirCall: Automated Conversational IVR in Airport Call Center Using Deep Transfer Learning

Soheil Tehranipour¹, Mohammad Manthouri^{2*}, Samaneh Yazdani³

1. MS Student, Department of Computer Engineering, North Tehran Branch, Islamic Azad University, Tehran, Iran. s.tehranipour@iau-tnb.ac.ir
2. Assistant Professor, Department of Electrical and Electronic Engineering, Shahed University, Tehran, Iran.
*Corresponding Author, mmanthouri@shahed.ac.ir
3. Assistant Professor, Department of Computer Engineering, North Tehran Branch, Islamic Azad University, Tehran, Iran. s.vazdani@iau-tnb.ac.edu

Abstract

Introduction: In this paper, we introduce the ParsAirCall toolkit, which is a tool for automatic recognition of Persian numbers in airport systems. It leverages deep transfer learning to improve performance in real and operational scenarios of voice-controlled smart telephone systems at airports across the country. In today's world, with the advancements in artificial intelligence, traditional systems for interacting with callers in telephone calls are not efficient, and this efficiency will be enhanced through automation and the automation of repetitive tasks.

Method: ParsAirCall distinguishes itself by surpassing competing models in the Persian language, achieving heightened accuracy with fewer parameters and optimized computing resources. Addressing the challenge posed by limited data for Persian speech recognition, we meticulously curated a 30-hour telephony dataset, serving as the cornerstone for training the final ParsAirCall model. Embracing the innovative QuartzNet architecture, our deep transfer learning strategy empowers ParsAirCall to capture nuanced features in Persian speech, ensuring superior performance in number recognition tasks associated with airport telephone calls.

Results: Experiments were conducted on both our collected telephony dataset and the Common Voice project, demonstrating ParsAirCall's efficiency in achieving a 2.7% WER (Word Error Rate) in number recognition in airport telephone calls.

Discussion: ParsAirCall emerges as a versatile tool, poised for seamless integration as a service into any Persian-language airport telephone system. Its practical application extends to number recognition in airport call centers, exemplifying the transformative impact of advanced technologies in streamlining communication processes within critical operational environments. ParsAirCall can be easily integrated as a service into any Persian-language airport telephone system, making it a practical tool for number recognition in airport call centers and telephone systems.

Keywords: Call center, automatic speech recognition, Deep transfer learning, Airport smart systems.



انجمن علمی تجارت الکترونیکی ایران

سامانه‌های پردازشی و ارتباطی چندرسانه‌ای هوشمند

Intelligent Multimedia Processing and Communication Systems (IMPCS)



واحد تهران

ParsAirCall: تلفن گویای محاوره‌ای خودکار در مرکز تماس فرودگاهی با استفاده از یادگیری انتقالی عمیق

دوره چهارم، تابستان ۱۴۰۲
شماره دوم، صص: ۱۱-۲۰

تاریخ دریافت: ۱۴۰۲/۰۱/۱۶
تاریخ پذیرش: ۱۴۰۲/۰۲/۲۷

سهیل تهرانپور^۱، محمد منثوری^{۲*}، سمانه یزدانی^۳

۱- کارشناسی ارشد، هوش مصنوعی و رباتیک، واحد تهران شمال، دانشگاه آزاد اسلامی، تهران، ایران. s.tehranipour@iau-tnb.ac.ir

۲- استادیار، گروه کنترل، دانشگاه شاهد، تهران، ایران. (نویسنده مسئول) mmanthouri@shahed.ac.ir

۳- استادیار، گروه کامپیوتر، واحد تهران شمال، دانشگاه آزاد اسلامی، تهران، ایران. s.vazdani@iau-tnb.ac.edu

چکیده: در دنیای امروز با عنایت به پیشرفت‌های هوش مصنوعی، سامانه‌های سنتی موجود در تعامل با مخاطبین در تماس‌های تلفنی کارآمد نخواهند بود و این ناکارآمدی با اتوماسیون و خودکارسازی فعالیت‌های تکراری بهبود خواهد یافت. در این مقاله، ابزاری با نام ParsAirCall معرفی شده است که ابزاری برای بازسازی خودکار گفتار فارسی اعداد در سامانه‌های فرودگاهی است که از یادگیری انتقالی عمیق برای بهبود عملکرد در محیط‌های واقعی و عملیاتی سامانه‌های هوشمند تلفنی گویا در فرودگاه‌های کشور می‌باشد. ParsAirCall با تعداد کمتر پارامترها و منابع محاسباتی، بهینه‌تر و نیز دقت بالاتر از مدل‌های رقیب عمل می‌کند. به دلیل محدودیت داده‌های موجود برای تشخیص گفتار در زبان فارسی، یک مجموعه دادگان ۳۰ ساعته تلفنی جمع‌آوری و برچسب‌گذاری شده و برای آموزش مدل نهایی مورد استفاده قرار گرفته است. آزمایش‌ها روی هر دو مجموعه داده تلفنی جمع‌آوری شده توسط ما و پروژه Common Voice انجام شدند، که کارایی ParsAirCall را در دستیابی به ۲,۷٪ میزان WER در تشخیص اعداد در تماس‌های تلفنی فرودگاهی نشان می‌دهد. ParsAirCall می‌تواند به راحتی به عنوان یک سرویس در سیستم‌های تلفنی فرودگاهی فارسی زبان ادغام شود، که آن را به یک ابزار قابل اجرا برای تشخیص اعداد در مراکز تماس و سیستم‌های تلفنی فرودگاهی تبدیل کند.

واژه‌های کلیدی: مرکز تماس، بازسازی خودکار گفتار، یادگیری عمیق انتقالی، سامانه‌های هوشمند فرودگاهی.

۱. مقدمه

ارموزه خدمات پس از فروش و پشتیبانی مشتریان یک بخش اساسی در اداره کسب و کارهای مدرن می‌باشد که به‌ویژه در سال‌های اخیر به دلیل مدل‌های جدید شغلی، استفاده گسترده از هوش مصنوعی، گسترش روند استفاده از تکنولوژی‌های نوظهور و بالاتر رفتن سطح انتظارات برای تجربه کاربری دگرگون شده است. [۱] ایده‌ها، پروژه‌ها و فناوری‌های جدید زیادی در حال کمک به رشد بازار هستند که همگی سعی در بهبود وضعیت موجود و رفع مانع می‌باشند. راه‌حل‌های سنتی مراکز تماس در عمل قادر به عملیات و پاسخگویی به مشتریان نخواهند بود؛ به‌عنوان نمونه می‌توان به افزایش فزاینده حجم تماس‌های تلفنی به مراکز تماس برای دریافت خدمات پشتیبانی اشاره نمود که مشخصاً پس از شیوع کووید-۱۹ و متعاقباً افزایش بی‌سابقه خدمات از راه دور، این صنعت را دچار آسیب جدی ساخت.

در پاسخ به وضعیت پیش‌آمده، بسیاری از کسب و کارها از فناوری‌های جدیدتر مانند VoIP و هوش مصنوعی برای اجرای عملیات خدمات پشتیبانی و ارائه سرویس به مشتریان خود استفاده می‌کنند. یکی از مهم‌ترین خدمات جدید در مراکز تماس مدرن، استفاده از تشخیص خودکار گفتار (ASR) است که به مدیران مراکز تماس کمک می‌کند تا انبوهی از داده‌های تلفنی از منابع مختلف مانند واکنش صوتی تعاملی (IVR) و امکان ضبط تماس‌ها را داشته‌باشند. به این ترتیب سامانه هوشمند نصب‌شده در مراکز تماس، به تماس‌گیرنده اجازه می‌دهد تا فقط از طریق گفتار خود با مرکز تماس ارتباط برقرار کند؛ این در حالی است که در گذشته و راهکارهای سنتی، از صفحه کلید برای وارد کردن داده‌های موردنظر استفاده می‌شد. در مراکز تماس نوین، مسيردهی و تکمیل تماس مخاطب، با تبدیل گفتار به متن و سپس پردازش آن متن برای ارجاع تماس به مقصد تماس‌گیرنده انجام می‌شود. علاوه بر این، در ترکیب با دیگر برنامه‌ها، فناوری‌های موجود به تماس‌گیرندگان اجازه می‌دهد تا خدماتی مانند بررسی موجودی حساب و احراز هویت خود را آنی و بدون نیاز به صحبت با اپراتور انسانی و یا ورود اطلاعات به‌صورت دستی انجام‌دهند که در نهایت منجر به تجربه بهتر مشتری و پشتیبانی ۷/۲۴ در مراکز تماس می‌شود. اقدامات مذکور نمونه‌ای از مواردی است که منجر به صرفه‌جویی در هزینه و افزایش سطح رضایت در خدمات مشتری خواهد شد. [۲]

از سوی دیگر، در چند سال اخیر، چندین روش تبدیل گفتار به متن، توسط مدل‌های End-to-End (E2E) با استفاده از شبکه‌های عصبی و تکنولوژی‌های جدید نوظهور ساخته شده‌اند. برای آموزش مدل‌های گفتار به متن از ابتدا، نیاز به مقدار قابل توجهی داده برچسب‌خورده و توان محاسباتی فوق‌العاده داریم؛ با وجود این که بیش

از ۶۰۰۰ زبان در سراسر جهان صحبت می‌شوند، تنها معدودی از زبان‌ها هستند که دارای مجموعه داده‌های غنی برای آموزش مدل‌های با کیفیت برای گفتار به متن می‌باشند. یادگیری انتقالی یک رویکرد معتبر برای توسعه مدل‌های گفتار به متن برای زبان‌های با منابع کم مانند فارسی است که از مدل آموزش داده‌شده برای یک زبان دیگر استفاده می‌کند. هدف ما از این رویکرد، جبران کمبود داده مورد نیاز به جهت آموزش مدل‌های گفتار به متن فارسی برای استفاده در مراکز تماس و مراکز تلفن می‌باشد. طبق تحقیقات صورت‌گرفته، گفتار به متن فارسی تا به امروز در مراکز تماس هنوز عملیاتی نشده‌است. با توجه به افزایش محبوبیت مراکز تماس در میان فارسی‌زبان، ما تأثیر یادگیری انتقالی بر سامانه گفتار به متن فارسی به‌خصوص چگونگی اتوماسیون کامل فرآیند تماس در سیستم پاسخ صوتی تعاملی در صورت برقراری ارتباط با تماس‌گیرنده برای ارتباط با مخاطبین را بررسی نموده‌ایم. در نهایت در این مقاله، نتایج نهایی یادگیری انتقالی در زبان فارسی را با در نظر گرفتن مدل‌های مختلف روز دنیا ارائه می‌دهیم. همچنین، ما کشف کردیم که اگر مجموعه داده نهایی کوچک‌تر، هم‌اندازه یا حتی به‌طور قابل‌ملاحظه‌ای بزرگتر از مجموعه داده آموزشی باشد، به‌طرز عجیبی، استفاده از یک مدل پیش‌آمورخته همچنان نسبت به آموزش کلیه داده‌ها از ابتدا، نتایج بهتری خواهد داد. نوآوری‌های مقاله حاضر به شرح زیر خواهد بود:

- ارائه یک سیستم تشخیص گفتار فارسی تک‌زبانه برای تشخیص اعداد فارسی با دقت بالا با استفاده از یادگیری انتقالی.
- عملکرد بهتر ParsAirCall در مقایسه با سایر معماری‌های توسعه داده‌شده برای زبان فارسی.
- با توجه به هدف ما جهت استفاده عملیاتی در مراکز تماس، ParsAirCall از نظر مصرف منابع سخت‌افزاری و بار پردازشی بهینه شده است.
- در راهکار نهایی، ParsAirCall جهت استفاده در مراکز تماس و هوش مصنوعی گفتگو محور به زبان فارسی پیاده‌سازی و بهینه شده است.

در بخش ۲ به بررسی مطالعات اخیر در حوزه بازشناسی گفتار خواهیم پرداخت و در بخش ۳ مجموعه دیتاست این پژوهش ارائه خواهد شد. بخش ۴ روش پیاده‌سازی و بخش ۵ تحلیل گزارش نهایی خواهد بود.

۲. پیشینه پژوهش

در طول دهه گذشته، پردازش تصویر، بازشناسی خودکار گفتار، پردازش زبان طبیعی و دیگر حوزه‌های هوش مصنوعی، پیشرفت عظیمی داشته‌است. این پیشرفت عمدتاً به دلیل حضور بلامنازع یادگیری عمیق و انواع مختلف آن می‌باشد. یادگیری عمیق، این امکان را می‌دهد که در صورت

ساخت مدل با داده‌های آموزشی متنوع و کیفیت بالا، پیش‌بینی‌ها بسیار دقیق صورت‌گیرد و همان‌طور که در پیشتر اشاره شد، با استفاده از راهکارهای ارائه‌شده، اغلب یکی از معماری‌های استاندارد و معروف موجود استفاده می‌شود. از سوی دیگر، اگر مجموعه داده بزرگی برای مسئله موردنظر در اختیار نباشد، چاره کار در استفاده از یادگیری انتقالی است؛ در این حالت، با استفاده از یک شبکه از پیش آموزش دیده بر روی حجم عظیمی از داده‌ها و تنظیم برخی از لایه‌های انتهایی آن، مسئله موردنظر قابلیت اجرا خواهد داشت.

اغلب تحقیقات اولیه با محوریت بازشناسی خودکار گفتار، در دهه‌های ۵۰ و ۶۰ میلادی بر روی کلمات به صورت ایزوله و وابسته به صحبت کاربر انجام شده است. تلاش‌هایی نیز در سطح واج صورت گرفت، اما با توجه به دشواری بالقوه این پروژه، موفقیت‌های چشمگیری حاصل نشد. به عنوان یکی از ابتدایی‌ترین موارد پیاده‌سازی شده، در اوایل سال ۱۹۵۲ در آزمایشگاه‌های موسوم به آزمایشگاه بل، سیستم‌های کلمات جداشده ارائه گردید که قابلیت تشخیص ارقام خوانده شده توسط یک مرد بالغ را داشت؛ نکته قابل توجه، نرخ دقت ۹۹٪ در تشخیص ارقام مجزای این روش می‌باشد. این سیستم برای تشخیص کل کلمات، تنها قادر به تشخیص در فرکانس‌های موجود بوده و هیچ الگویی از هجاها یا صامت‌ها یا مصوت‌ها یا هر نوع واحد زیرکلمه‌ای در این سیستم‌ها وجود ندارد. کلمه به عنوان یک واحد در نظر گرفته می‌شود و در طول طبقه‌بندی، تمام کلمات برای یافتن بهترین تطابق با یکدیگر مقایسه می‌شوند. [۳]

مدل‌های مارکوف پنهان (HMM) در دهه ۷۰ میلادی، مدل‌های مخلوط گاوسی (GMM) نیز در دهه ۸۰ میلادی معرفی شدند، اما تا دهه ۹۰ میلادی زمان لازم بود تا هر سه الگوریتم در یک راهکار ادغام شوند؛ که آن را HMM Toolkit و به اختصار HTK نامیدند. HTK، اولین سامانه جهت ترکیب تمام اجزای اصلی شناسایی گفتار مبتنی بر GMM-HMM مدرن بود که در سال ۱۹۹۲ به عنوان الگوریتم متن‌باز در اختیار عموم قرار گرفت. [۴] استفاده از روش مرسوم تشخیص گفتار، یعنی HMM به همراه GMM، به حدی رسید که بر روی پروژه‌های تشخیص گفتار پیوسته مستقل از گوینده با واژگان بسیار گسترده کارایی بالایی پیدا کرد و عمده پروژه‌های تشخیص خودکار گفتار، از این الگوریتم استفاده نمی‌کرد. اما در این گونه روش‌ها فرض‌های مختلفی در نظر گرفته می‌شود که برخی غیرواقعی بوده و در عمل چالش بزرگی را پیش‌روی کاربردهای تجاری و صنعتی می‌گذارد. به عنوان مثال فرض می‌شود که ویژگی‌های استخراجی در یک بخش الفبایی با یکدیگر همبسته نیستند که این مورد نادرست است. از طرف دیگر، به منظور پیاده‌سازی عملیاتی پیچیدگی‌های بسیار زیادی را برای مدل‌هایی به همراه خواهند داشت. روش‌های مبتنی بر شبکه‌های عصبی راهی برای کاهش وابستگی‌های سیستم به فرضیات غیرواقعی درباره گفتار معرفی نمودند. زیرا شبکه‌های عصبی، توانایی یادگیری همبستگی بین ویژگی‌ها را دارا می‌باشد؛ بنابراین گامی رو به جلو برداشته شد و کارایی بیشتر روش HMM موجب تولید سیستم‌های تشخیص گفتار بر پایه آن شد.

اما بازشناسی گفتار، پیوسته یک موضوع حل نشده باقی مانده است که همچنان چالش‌های بی‌شماری در تولید یک سامانه کامل با دقت بالا را پیش‌روی خود می‌بیند. مشکلات موجود در دنیای واقعی عبارتند از: اندازه بزرگ دیکشنری، تشخیص‌دهنده‌های مستقل از گوینده (در مقابل وابسته به گوینده که ساده‌تر است)، پدیده گفتار ناخواسته (کلمات خارج از دیکشنری) لهجه و گویش‌های ناآشنا، حالت‌های مختلف صحبت (لحن)، نویز پس‌زمینه، صداهای انعکاسی اتاق، صحبت‌های عامیانه که در دایره لغات آموزشی نیستند، همگی بخشی از این چالش‌ها به-شمار می‌روند. [۵]

شبکه‌های عصبی مصنوعی (ANN) را می‌توان برای طبقه‌بندی واحدهای گفتاری مانند کلمات یا واج‌ها به کار برد اما نمی‌توان از آن‌ها برای تشخیص گفتار استفاده کرد، زیرا تعداد دنباله‌های کلمات ممکن در یک گفتار نامحدود است و شبکه عصبی دارای محدودیت اندازه است؛ حال آنکه افزایش اندازه ورودی در یک شبکه عصبی موجب افزایش تعداد نورون‌های لایه مخفی و پارامترها می‌شود که چالش‌ها و سختی‌های خود را دارد. ایده اصلی این گونه است که از شبکه‌های عصبی برای طبقه‌بندی‌های محلی و از HMM برای مدل کردن ساختار دنباله‌ای داده‌ها استفاده شود. HMM توانایی بخش‌بندی خودکار دنباله‌های ورودی در طول آموزش را دارد و حتی می‌تواند یک روش اصولی برای تبدیل کلاس‌بندی‌های شبکه به دنباله‌ای از برچسب‌ها تهیه کند. بنابراین برخلاف شبکه‌های عصبی، ترکیب‌های ANN-HMM را می‌توان به-صورت مستقیم در کارهای طبقه‌بندی زمانی مانند بازشناسی خودکار گفتار به کار برد. [۶]

در ادامه روش‌های پیشین، انتشار مقاله هینتون و همکاران در سال ۲۰۱۲ با عنوان «شبکه‌های عصبی عمیق برای شناسایی گفتار»، قدرت شبکه‌های عصبی به‌ویژه یادگیری عمیق را در تشخیص گفتار تثبیت نمود. [۷] نخستین رویکردها برای رسیدن به این هدف، سعی در درک معنای گفتار با استفاده از تبدیل آن به متن و سپس بررسی قواعد نوشته شده توسط انسان بود. با این حال، چنین سیستم‌های قاعده‌محوری، محدود به حوزه‌های خاصی بودند که برای آن‌ها طراحی شده بود و عموماً قادر به مقابله با ورودی‌های غیرمنتظره یا ورودی‌های جدید نبودند و در نهایت ثابت شد که برای به دست آوردن پیچیدگی زبان طبیعی بسیار محدودند. به کل بخش استخراج ویژگی و محاسبه احتمال برچسب‌ها پردازش صوتی و بخش محاسبه احتمال (انتخاب برچسب بر حسب ویژگی‌های استخراجی)، مدل صوتی گفته می‌شود؛ که در [۷] از آن استفاده شده است.

از طرف دیگر، طی ۲۰ سال گذشته، رویکرد آماری بازشناسی گفتار خودکار نیز متداول شده است که از مدل‌های ریاضی برای یادگیری خودکار قوانین استفاده می‌کند. در نتیجه، به جای نوشتن قوانین سخت و پیچیده که انعطاف زیادی هم ندارند، فرآیند پیشرفت بازشناسی گفتار خودکار، به سمت ایجاد ویژگی‌هایی که بیانگر روابط بین داده‌هاست، هدایت گردید. با این همه، تهیه و استخراج ویژگی‌های مهندسی زمان‌بر

است، زیرا این ویژگی‌ها معمولاً وظیفه‌محور بوده و به تخصص در حوزه مورد نظر نیاز دارند.

یک فرض اساسی در یادگیری آماری باناظر این است که داده‌های آموزشی برچسب‌دار دارای توزیع مستقل و یکنواخت هستند که در نهایت مدل را بر روی نمونه‌های آزمایشی، اجرامی کنند. اما در صورت پایین بودن خطای آموزش، نمی‌توان عملکرد مناسب در مواجهه با داده‌های آزمون را تضمین نمود. به‌ویژه زمانی که با داده‌های صوتی و متنی سروکار داریم، این فرض معمولاً به این معنی است که داده‌های آموزشی و آزمایشی باید از یک دامنه باشند. مثلاً، اگر کسی بخواهد یک اپلیکیشن با کاربرد تشخیص موجودیت‌های اسمی (NER) برای وبلاگ‌های شخصی ایجاد کند، در حالت عادی باید به‌جای آموزش مدل بر روی مقالات خبری مالی (که اغلب برچسب‌دار هستند)، آن را با استفاده از موجودیت‌های اسمی آموزش دهد زیرا توزیع داده‌ها در حوزه خبری و حوزه‌های متفاوت است. [۸] یادگیری عمیق باناظر، متکی بر مقدار کافی از داده برچسب‌گذاری شده برای آموزش است؛ به همین دلیل، زبان‌شناسان محاسباتی تلاش زیادی کرده‌اند تا پیکره‌های متنی برچسب‌گذاری شده مختلفی برای مقاصد گوناگون تهیه کنند. [۹] از سوی دیگر، پیکره‌های ثبت‌شده موجود، حتی در صورت کار بر روی دامنه‌های دیگر نیز، مفید خواهند بود. بنابراین، استفاده از روش‌هایی که بتوانند به‌عنوان میان‌بر، تولید نرم‌افزارهای نهایی را سرعت دهند، می‌تواند ارزش افزوده بالایی را به ارمغان آورد. بنابراین، برای رسیدگی به مشکل گلوگاهی با فقدان وجود تعداد کافی از داده‌های برچسب‌دار، یک سوال تحقیقاتی مهم و جالب این است که چگونه از داده‌های برچسب‌دار موجود از دامنه‌های مرتبط برای آموزش طبقه‌بندها در دامنه‌ای که ما به آن علاقه داریم، استفاده کنیم. به‌طور واضح‌تر، دامنه‌ای که داده‌های برچسب‌دار کافی ندارد ولی مورد علاقه ماست را هدف و دامنه‌ای که داده‌های برچسب‌دار کافی دارد، دامنه منبع تعیین می‌کنیم. جهت حل‌اولین روش‌های شناسایی گفتار پیوسته، از تابع هدف طبقه‌بندی زمانی متصل یا CTC استفاده کردند [۱۱]. مدل توالی به توالی با رویکرد مکانیزم توجه (Attention Mechanism) خیلی زود تبدیل به مدل محبوب محققین شد. یک نسخه قابل توجه از مدل توالی به توالی با توجه، LAS است که در ادامه پس از استفاده از در مدل‌های بازگشتی، روش Wav2Letter معرفی و بر روی لایه‌های پیچشی و اصلاح تابع اتلاف CTC ساخته شد. [۱۲] در طول پنج سال گذشته، مدل‌های پیوسته‌تر که کارآمدتر شده‌اند، اما معمولاً در مجموعه داده‌های عظیم آموزش می‌بینند که تنها در دسترس شرکت‌های بزرگ هستند. لذا، پژوهش حاضر با هدف چگونگی تعمیم این مدل‌های پیوسته، برای مجموعه داده‌های کوچک‌تر، انجام شده است. در اواخر سال ۲۰۱۴، روش DeepSpeech تحت نظر اندرو انگ، پیشنهاد شد که از مدل شبکه عصبی بازگشتی (RNN) استفاده نموده است. [۱۳]

مشکل مذکور، راهکار موسوم به تطبیق دامنه پیشنهاد شده است که از داده‌های برچسب‌دار دامنه منبع برای کمک به آموزش طبقه‌بند در دامنه

هدف استفاده می‌کند. [۱۰] این موضوع یکی از زیر مجموعه‌های اصلی یادگیری انتقالی است. طبق مقالات و تحقیقات منتشر شده در حوزه یادگیری ماشین به‌ویژه در سال‌های اخیر، شاهد رشد روزافزون توجه‌ها به یادگیری انتقالی و تطبیق دامنه هستیم. در محدوده پنج سال گذشته، شبکه‌های عصبی عمیق، به‌عنوان یک دسته خاص از مدل‌های یادگیری ماشین، با توجه به کارایی فوق‌العاده‌ای که از خود نشان دادند به مدل محبوب و منتخب محققین و متخصصین علوم داده و هوش مصنوعی تبدیل شدند که با توجه به نوع عملکردشان، نیاز به مهندسی ویژگی و تعیین ویژگی‌ها به‌صورت دستی و توسط نیروی انسانی خبره را کاهش می‌دهد. در نتیجه انرژی نیروی انسانی بر تعیین مناسب‌ترین معماری و تنظیمات آموزشی برای هر کار متمرکز خواهد شد. نسل فعلی مدل‌های بازشناسی گفتار خودکار مبتنی بر شبکه عصبی این برتری را دارند که می‌توانند بر روی مقادیر بسیار زیاد داده‌های برچسب‌دار آموزش ببینند. در بازشناسی گفتار خودکار و متعاقب آن، پردازش زبان طبیعی و دیگر حوزه‌های یادگیری ماشین، روش استاندارد برای آموزش یک مدل به این ترتیب خواهد بود که ابتدا باید تعداد زیادی از نمونه‌های جمع‌آوری شده برچسب‌گذاری شود. سپس این داده‌ها به‌عنوان ورودی به یکی از مدل‌های موجود ارائه می‌شود و این مدل طی فرآیند آموزش، یاد می‌گیرد که از ورودی به خروجی نگاشت شود؛ این مورد با عنوان یادگیری باناظر شناخته می‌شود. تعداد نمونه‌هایی که باید برای هر موضوع جدید برچسب زده شوند، با توجه به فراوانی زبان‌ها، وظایف و حوزه‌ها در دنیای واقعی و تنظیمات مدل پیشنهادی، متفاوت بوده و به‌عنوان چالش‌های دنیای واقعی تلقی می‌شوند. در سال‌های اخیر، بازشناسی خودکار گفتار پیوسته بر روش‌های سنتی، ترجیح داده شده است زیرا برای تخمین پارامترها تنها یک تابع جهت محاسبه Loss وجود دارد. این مقدار Loss، به ما اجازه می‌دهد تا مستقیماً مسئله مورد نظرمان را بدین ترتیب، مدل کنیم و بهترین رونوشت ممکن برای یک فایل صوتی را به دست بیاوریم. با این حال، با حذف HMM ها از رویکرد ترکیبی، مدل‌های پیوسته، به اطلاعات هم‌ترازی که معمولاً در آموزش مدل‌های سنتی مورد استفاده قرار می‌گیرند، دسترسی ندارند. به این ترتیب، مدل تشخیص گفتار سرتاسری، باید هم‌ترازی متن و صدا را در طول آموزش بیابد؛ که قطعاً کار ساده‌ای نیست.

در ادامه روش مذکور، مدل توسعه‌یافته آن با نام DeepSpeech2 در سال ۲۰۱۶ ارائه گردید که نسبت به روش پیشین تا ۴۳٪ دقیق‌تر و ۷ برابر سریع‌تر می‌باشد. این مدل بر روی ۱۱,۹۴۰ ساعت داده زبان انگلیسی و ۹,۴۰۰ ساعت داده زبان چینی آموزش داده شده است. یکی دیگر از ابداعات این روش استفاده از تابع بهینه‌ساز جدیدی به نام SortaGrad است که از استراتژی آموزش برنامه‌درسی بهره می‌گیرد. معماری اصلی این مدل، در مجموع شامل ۱۱ لایه است

از فایل‌های صوتی شامل گفتار مرجع و متن متناظر آن گفتار. با توجه به نیاز الگوریتم‌های یادگیری عمیق به دادگان آموزشی زیاد و رشد قدرت محاسبات توسط کارت‌های گرافیکی در سال‌های اخیر، جهت به-دست‌آوردن مدلی مطلوب، حجم دادگان آموزشی بسیار مهم است. مجموعه دادگان گفتاری جامعی تهیه و جمع‌آوری شده‌است که شامل ۳۰ ساعت دادگان گفتاری برچسب‌دار در زبان فارسی می‌باشد. همان‌طور که پیشتر تبیین گردید، یکی از چالش‌های موجود، نبود داده‌های کافی در زبان فارسی جهت آموزش و ساخت مدل برای اجرا در این زبان است؛ بنابراین ارتباط تنگاتنگ با صنعت می‌تواند بخش زیادی از این مسئله را هموارسازد. با عنایت به توافقات صورت‌گرفته با شرکت‌های دانش‌بنیان فعال در حوزه راهکارهای نسل جدید مخابراتی و مراکز تماس، امکان تأمین بخشی از این داده‌ها جهت ساخت یک مدل جدید در زبان فارسی حاصل شد. طبق جدول ۱، ۳۰ ساعت فایل صوتی ضبط‌شده در مکالمات تلفنی مراکز تماس فارسی زبان، شامل ۱۰ هزار فایل صوتی (هر فایل صوتی با مدت زمان تقریباً ۱۰ ثانیه) گردآوری و برچسب‌گذاری شده-است. این مجموعه دادگان شامل فایل‌های صوتی تلفنی گسترده‌ای از اعداد ارقام به‌صورت مجزا، شماره تلفن، کد ملی، شماره شبا و ترکیب‌های مختلف عددی در زبان فارسی است.

جدول ۱: مجموعه دادگان

مدت زمان فایل صوتی (ساعت)	تعداد فایل صوتی	دیتاست مورد استفاده
۲۵۵	۳۷۶۵	Common Voice
۳۰	۱۰۰۰۰	مرکز تماس، برچسب‌گذاری شده
۳۸۵	۱۲۶۷۵	مجموع

یکی دیگر از منابع دادگان ما در این تحقیق، پروژه متن‌باز Common Voice است که توسط شرکت Mozilla برای ایجاد یک پایگاه داده رایگان جهت تشخیص گفتار در سال ۲۰۱۷ آغاز شده-است. [۲۰] این پروژه توسط داوطلبانی پشتیبانی می‌شود که نمونه جملات را با میکروفون ضبط می‌کنند و فایل‌های ضبط‌شده سایر کاربران را بررسی و تأیید می‌کنند. جملات رونویسی شده در یک پایگاه داده صوتی که تحت مجوز مالکیت عمومی CC0 موجود است جمع‌آوری خواهد شد. این مجوز اطمینان می‌دهد که توسعه‌دهندگان می‌توانند از پایگاه داده برای برنامه‌های صوتی به متن بدون محدودیت و هزینه استفاده کنند. تمامی ۹،۲۸۳ ساعت ثبت‌شده در این مجموعه داده نیز شامل فراداده جمعیتی مانند سن، جنس و لهجه است که می‌تواند به آموزش دقت موتورهای تشخیص گفتار کمک کند. این مجموعه داده در حال حاضر، شامل ۷،۳۳۵ ساعت تأییدشده به ۶۰ زبان است، اما همچنان در حال گسترش و افزودن داده‌های صوتی می‌باشد.

که در ابتدا سه لایه پیچشی، سپس هفت لایه بازگشتی دوطرفه و در نهایت یک لایه تمام متصل را شامل می‌شود. خروجی نهایی همانند مدل DeepSpeech به الگوریتم CTC تحویل داده می‌شود. [۱۴]

هیسائو و همکاران در سال ۲۰۲۰، در ادامه روش LAS، نسخه آنلاین آن را ارائه نمودند تا نگاه‌ها به نسخه‌های بعدی این روش یعنی DeepSpeech3 دوخته شود که طبق گفته شرکت بزرگ Baidu می‌تواند گام بزرگ دیگری در بازشناسی خودکار گفتار در تمامی زبان‌ها باشد. [۱۵]

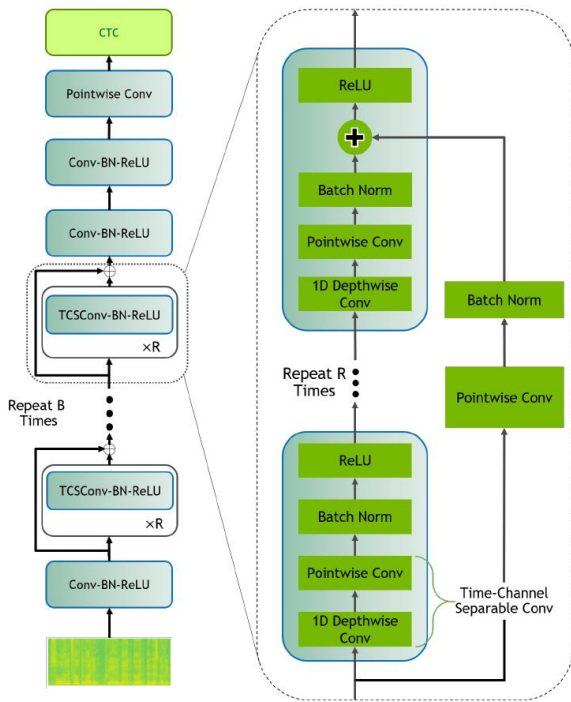
روش دیگری نیز توسط دپارتمان هوش مصنوعی شرکت فیس‌بوک، ارائه و با نام تجاری Wav2letter معرفی شد. کولوبرت و همکاران در اواخر سال ۲۰۱۶ آن را ارائه کردند [۱۶] که با تغییر نگرش در معماری و نحوه استفاده از بازنمایش داده‌های صوتی، توانست در پیچه‌ای برای شروع مطالعات بعدی بگشاید. این روش دو سال بعد یعنی سال ۲۰۱۸ در مقاله‌ای با نام Wav2letter++ ادامه یافت که خود را با عنوان سریع‌ترین روش متن‌باز معرفی نمود. طبق اعلام تیم فنی فیس‌بوک، روش Wav2Vec2.0 تنها با داشتن ۱۰ دقیقه داده صوتی برچسب‌دار می‌تواند به دقت بالایی در بازشناسی خودکار گفتار برسد. [۱۷] در این مدل‌ها که عموماً در یادگیری انتقالی، به‌ویژه برای زبان‌های با داده برچسب‌دار کم استفاده خواهند شد، امکان بهبود عملکرد سامانه‌های بازشناسی در داده‌های تلفنی نیز خواهد داشت [۱۸]

اخیراً نیز الگوریتم Whisper توسعه‌یافته توسط OpenAI نتایج چشم‌گیری را در تشخیص گفتار با استفاده از یادگیری خودنظارتی به همراه یک تکنیک خوشه‌بندی نوآورانه نشان داده‌است. هنگام مقیاس‌بندی به ۶۸۰،۰۰۰ ساعت از نظارت چندزبانه و چندوظیفه‌ای، مدل‌های حاصل تعمیم خوبی به بنچمارک‌های استاندارد دارند و در مقایسه با انسان‌ها، مدل‌ها عملکرد دقیقی داشته‌اند. این پیشرفت‌ها در زمینه یادگیری خودنظارتی قابلیت توسعه بسیار بزرگی را برای پروژه‌های علمی که از یادگیری ماشینی استفاده می‌کنند، دارند زیرا اجازه می‌دهند مدل‌های دقیق‌تری با داده‌های کمتر برچسب‌خورده آموزش داده-شوند. [۱۹]

با وجود مطالعات اخیر در زمینه تشخیص گفتار فارسی، چالش‌ها همچنان در دستیابی به تشخیص به‌صورت زمان واقعی باقی‌مانده‌اند، به‌خصوص در مراکز تماس که داده‌های تلفنی ممکن است چالش‌های بزرگی داشته باشند. به عبارت دیگر، نیاز به تحقیقات بیشتری برای توسعه سیستم‌های تشخیص گفتار اتوماتیک واقعی و تعاملی در محیط‌های عملیاتی مراکز تماس و همچنین مدیریت پیچیدگی‌های داده‌های تلفنی وجود دارد. با حل این چالش‌ها، سیستم‌های تشخیص گفتار خودکار، پتانسیل تحولی در مراکز تماس به زبان فارسی را دارند و می‌توانند ارتباط مشتریان و عوامل مرکز تماس را بهبود بخشند.

۳. مجموعه دادگان

داده مورد نیاز جهت بررسی و پیاده‌سازی و آموزش یک سیستم بازشناسی گفتار دارای ساختار مشخصی است به شکل تعداد مشخصی



شکل ۱: معماری مدل QuartzNet [۲۱]

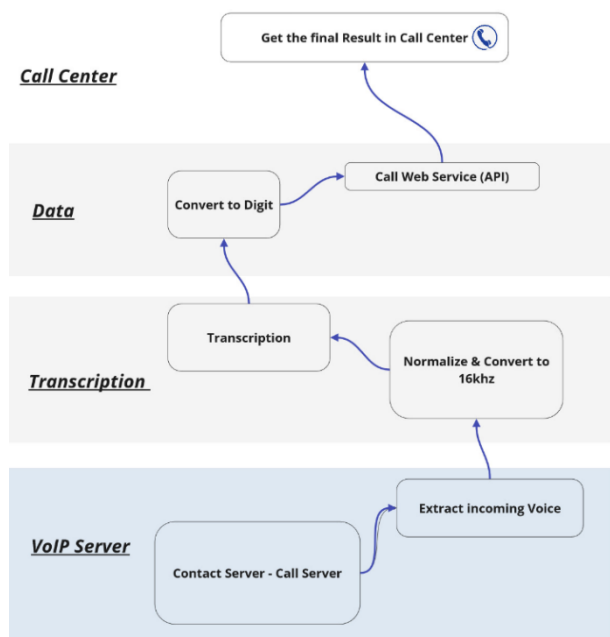
۴. پیاده‌سازی و ساخت مدل

همان‌طور که در تمامی زوایای زندگی همه ما مشهود است، استفاده از سیستم‌های کامپیوتری و نیز وسایل شخصی هوشمند به‌طور فزاینده‌ای رواج یافته‌اند. همچنین، به‌دلیل بسیاری از کاربردهای بالقوه این سیستم‌ها، یکی از جنبه‌های اصلی ارتباطی افراد که مکالمه و گفتگو می‌باشد نیز پایه‌ای بقیه کاربردها، باید از علوم روز و صدا البته هوش مصنوعی، استفاده نماید. هر نرم‌افزاری با محوریت کاربردی مکالمه و گفتگو، از چندین ماژول اساسی دیگر مانند تشخیص خودکار گفتار تشکیل شده است که مدل‌های استفاده‌شده هر چقدر سبک‌تر و بهینه‌تر باشند، نتیجه نهایی سریع‌تر و جذاب‌تر خواهد بود؛ به‌ویژه در جایی که اکثر دستگاه‌ها کوچک و کم‌حافظه با قدرت پردازش محدود هستند. با این حال، بیشتر مدل‌های پیشرفته بسیار بزرگ و سنگین هستند که معمولاً دارای چند صد میلیون پارامترند. با توجه به محدودیت‌های فعلی دستگاه‌ها، عملاً استقرار آن‌ها در مقیاس بزرگ دشوار و حتی در موارد زیادی، غیرممکن خواهد بود. برای مقابله با چالش فوق، از الگوریتم QuartzNet استفاده خواهیم کرد که یک معماری مبتنی بر یادگیری عمیق جهت پروژه بازشناسی خودکار گفتار می‌باشد. [۲۱] این مدل یک مدل End-to-End است که بزرگترین مدل آن، تقریباً ۱۹ میلیون پارامتر دارد و بادقت بسیار قابل‌قبولی، در مجموعه دادگان LibriSpeech، به ۲٫۶٪ خطای کلمه (WER) با استفاده از الگوریتم QuartzNet پردازش زبان طبیعی Transformer-XL می‌رسد. مدل QuartzNet همچنین توانایی یادگیری انتقال خوبی را نشان داده است و می‌تواند از

در این تحقیق برای زبان فارسی از ورژن ۱۳ پروژه Common Voice شامل ۳۵۵ ساعت داده صوتی تأییدشده و در مجموع ۳۹۲ ساعت داده صوتی، استفاده شده که به‌صورت رایگان و متن‌باز در اختیار همگان قرار گرفته است. هر ورودی در مجموعه داده شامل یک فایل MP3 و متن مربوط به آن است.

علاوه بر این، پروژه Common Voice در زبان فارسی فاقد یک استاندارد و راهنمای کامل جهت برچسب‌گذاری است و این مسئله موجب مشکلات قابل‌توجهی در فرآیند آموزش مجموعه داده می‌شود. مسئله نرمال‌سازی متون فارسی از بسیاری جهات حائز اهمیت است. متاسفانه به‌دلیل شرایط زبان فارسی که در ساخت کلمات جدید وابستگی زیادی به پیشوندها و پسوندها دارد نحوه نوشتاری کلمات تنوع زیادی دارد. به‌عنوان مثال پسوند جمع "ها" به سه صورت در متون وجود دارد که این تنوع در لغت‌نامه باعث می‌شود یک کلمه به سه صورت نوشته شود و دقت سیستم را کاهش دهد. مثلاً جمع کلمه "درخت" به سه صورت "درخت‌ها"، "درخت ها" و "درختها" نوشته می‌شود که در نگارش زبان فارسی تنها صورت درست با نیم‌فاصله است. یکی دیگر از مشکلات پیش روی پردازش متن در زبان فارسی، تنوع در یونیکدهای زبان فارسی ناشی از تفاوت در کیبوردهای مختلف و مشابهت حروف عربی و فارسی در این کیبوردهاست. حروفی مانند "ی" و "ک" بیشترین گوناگونی را از لحاظ یونیکد دارند. این نرمال‌سازی متون فارسی مستقیماً در دقت و نتایج لغت‌نامه آوایی و مدل زبانی تأثیری گذارد و در نهایت دقت سیستم بازشناسی گفتار را افزایش خواهد داد. بسیاری از نمونه داده‌های صوتی در مجموعه داده مذکور نیاز به نرمال‌سازی متون دارد که شامل حذف ویرگول‌ها و علائم نگارشی، علامت تعجب و کاراکترهای غیرواژه‌نامه‌ای دارد. همچنین تبدیل نمادهای انگلیسی و عربی به فارسی، حذف ایموجی‌ها، جایگزین کردن متن‌های فارسی قابل-فهم با نمادهای ریاضی، تبدیل تاریخ و زمان به متن، جایگزینی اختصارات با تلفظ آن‌ها و حذف برخی از کاراکترهای عربی در مجموعه داده مانند تنوین از اقدامات صورت‌گرفته جهت نهایی‌سازی مجموعه داده فوق است. ماهیت داده صوتی مناسب آموزش سیستم‌های بازشناسی خودکار گفتار نمی‌باشد.

تنظیم نرخ نمونه‌برداری: نرخ نمونه‌برداری باید میان تمامی داده‌ها یکسان و برابر عدد ۱۶ هزار باشد.



شکل ۲: فلوجارت خروجی مدل ParsAirCall از ابتدا تا انتها

تنظیم تعداد کانال: باید تمامی داده‌های ورودی دارای یک کانال صوتی باشند که این مهم به‌ویژه در مکالمات مراکز تماس از اهمیت زیادی برخوردار است.

آموزش مدل شبکه عصبی عمیق: در این بخش پس از پیش‌پردازش داده‌ها و ایجاد یک مدل پایه، مدل‌های شبکه عصبی عمیق پیاده‌سازی شده و درصد خطا بر روی یک داده آزمایشی ثابت بدست آمده است. این مجموعه دادگان آزمایشی شامل ۱۰۰۰ فایل صوتی، به‌صورت کاملاً تصادفی از بین داده‌های جمع‌آوری شده تلفنی مراکز تماس انتخاب شده است. ابزارهای مورد استفاده جهت پیاده‌سازی مدل‌های شبکه عصبی عمیق و آموزش این مدل‌ها در ادامه ارائه خواهد شد. به‌منظور اجرای پروژه آموزش مدل‌ها نیز از دو عدد GPU مدل Nvidia RTX3070 استفاده شده است و نحوه آموزش مدل به صورت موازی و Multi-GPU با استفاده از کتابخانه PyTorch Lightning و قابلیت DDP (Distributed Data Parallel) صورت گرفته است. همچنین ابزارهای دیگری نیز جهت پیش‌پردازش داده و پیاده‌سازی مدل‌های مبتنی بر یادگیری عمیق در این تحقیق به‌کاررفته‌اند که فهرست آن در جدول ۳ آورده شده است.

استفاده از مدل‌های زبانی جهت تبدیل حروف به ارقام جهت استفاده در مراکز تماس: در فاز انتهایی این تحقیق، به ایجاد زیرساختی جهت ارائه خروجی نهایی به‌صورت ارقام موردنظر جهت

یک نمونه از پیش‌ذخیره‌شده در یک مجموعه داده کوچک‌تر با دامنه محدود و محتوای تخصصی‌تر تنظیم‌شود تا در نهایت به دقت بسیار بالا در تشخیص خودکار گفتار دست‌یابد. QuartzNet از معماری Jasper مشتق شده است، هر دو مدل، مبتنی بر پیچش هستند که از بلوک‌های پیچشی، نرمال‌سازی دسته‌ای، ReLU و Dropout تشکیل شده‌اند و در انتها یک لایه با مقدار هزینه ناشی از CTC دارند.

جدول ۲: نمونه‌ای از داده‌های نرمال‌سازی شده

برچسب موجود در دیتاست اصلی COMMON VOICE	برچسب نرمال‌سازی شده جهت آموزش PARSAIRCALL
آنقدر خسته‌ام که ممکن است بیفتم.	آنقدر خسته ام که ممکن است بیفتم.
بسیار مشتاق دیدار ملکه‌ی زیبایی بودم!	بسیار مشتاق دیدار ملکه زیبایی بودم
عمداً آهسته از پله‌ها پایین رفت	عمدا آهسته از پله ها پایین رفت
صادقانه نظریه خود را اظهار کردیم	صادقانه نظریه خود را اظهار کردیم
مراتب تسلیم صمیمانه‌ی مرا بپذیرید	مراتب تسلیم صمیمانه مرا بپذیرید
اسطوره‌ی رستم و سهراب	اسطوره رستم و سهراب
برنامه‌ی جبران عقب‌افتادگی برای کودکان کوششین	برنامه جبران عقب افتادگی برای کودکان کوخ نشین
چهره‌ی چروکیده‌ی یک پیرمرد	چهره چروکیده یک پیرمرد
اختلاف عیقده‌ای، در بین ما نیست	اختلاف عیقده ای در بین ما نیست
او اهمیت آموزش را مورد تأکید قرار داد.	او اهمیت آموزش را مورد تاکید قرار داد.
زیرساخت غمانگیز این شعر	زیرساخت غم انگیز این شعر
جلسه‌ی کمیته ساعت چهار خاتمه یافت.	جلسه کمیته ساعت چهار خاتمه یافت.
به شدت باد تندی می‌آمد	به شدت باد تندی می آمد
یکی از جاسوسان کاگ ب	یکی از جاسوسان کا گه به
نهد و هیجده شصت و شیش	نهد و هجده شصت و شش

سیگنال خام گفتار که معمولاً با نرخ فرکانس ۱۶ کیلوهرتز (میکروفون موبایل یا کامپیوتر) یا ۸ کیلوهرتز (سیستم‌های تلفنی) نمونه‌برداری شده است، پس از قاب‌بندی و پنجره‌گذاری که معمولاً با نرخ ۲۵ میلی‌ثانیه و میزان همپوشانی ۱۰ تا ۱۵ میلی‌ثانیه انجام می‌شود را به‌عنوان ورودی به‌کار می‌بریم. سپس از قاب‌ها ویژگی‌های آن را استخراج می‌کنیم. به‌طور معمول از ۱۲ ضریب اول MFCC، مشتق اول و دوم آن به‌عنوان ویژگی استخراجی استفاده می‌نماییم و آن را به بخش بعدی انتقال می‌دهیم تا احتمال برچسب آن بخش از گفتار را محاسبه کند و برچسب مورد نظر را خروجی دهد.

نام مدل	تعداد پارامتر (میلیون)	زمان خروجی (ثانیه)	سایز مدل (MB)	WER (%)
ParsAirCall	۱۹	۱	۶۷	۲,۷
Whisper - Tiny	۳۹	۲,۲	۷۲	۴,۴
Whisper - Medium	۷۶۹	۲,۷	۱۴۲۰	۶,۱
Wav2vec 2.0	۹۱	۲	۳۶۳	۹,۵
DeepSpeech 2	۵۲	۲	۴۹۵	۱۵,۹

طبق جدول ۴، روش پیشنهادی پس از پیاده‌سازی یادگیری انتقالی و Fine-Tuning از سایر روش‌های موجود، بهینه‌تر عمل نموده و با در نظر گرفتن مدت زمان لازم جهت ارائه خروجی، برای استفاده در مرکز تماس عملیاتی فرودگاه‌های کشور پیشنهاد می‌گردد. همچنین لازم به ذکر است که میزان منابع مصرفی توسط مدل پیشنهادی، با توجه به سایز کمتر مدل و تعداد پارامترهای به‌مراتب کمتر، قابلیت اجرا در محیط‌های واقعی را داراست به‌گونه‌ای که در هسته تلفنی مرکز تماس تعبیه‌گردد. برای زبان فارسی، حیاتی است که مجموعه داده به‌خوبی مشخص شده‌ای وجود داشته باشد که این عوامل را در نظر بگیرد. سیستم تشخیص گفتار فارسی تولید شده، باید برای شناسایی انواع لهجه‌ها، گویش‌ها، و واژگان محاوره‌ای خاص در زبان فارسی آموزش داده شود تا WER کمتری به دست آید. ارزیابی باید به‌طور ایده‌آل با استفاده از چندین مجموعه آزمون انجام شود که شامل یک مجموعه گسترده از گویندگان، گویش‌ها، و شرایط گفتاری است تا درک جامعی از کارایی سیستم به دست آید.

۶. نتیجه‌گیری و جمع‌بندی

در این مقاله، ابزار ParsAirCall به‌عنوان یک ابزار تشخیص هوشمند اعداد فارسی در سامانه‌های فرودگاهی ارائه شده است که با استفاده از یادگیری انتقالی عمیق ساخته شده است. ParsAirCall با توجه به کمترین تعداد پارامترها و نیاز به منابع محاسباتی کمتر و نیز دقت بالاتر از مدل‌های رقیب، از مدل‌های مناسب برای محیط‌های عملیاتی است که می‌تواند به راحتی به‌عنوان یک سرویس در سیستم‌های تلفنی فرودگاهی فارسی زبان ادغام شود. ParsAirCall به دلیل Fine-Tuning بر روی ۳۰ ساعت داده تلفنی جمع‌آوری شده و نیز داده‌های پروژه Common Voice پس از نرمال‌سازی آن به‌مراتب به دقت بالاتری بر روی داده‌های واقعی در فاز آزمایشی رسیده است. همچنین بهره‌مندی از ایده یادگیری انتقالی و استفاده از مدل Pre-Trained مناسب، دلیل دیگر برای دستیابی به مدلی بهینه است. از ParsAirCall می‌توان به‌عنوان ابزاری برای تشخیص اعداد فارسی در مراکز تماس و سیستم‌های تلفنی فرودگاهی استفاده کرد. به‌منظور توسعه پژوهش‌های صورت گرفته در حوزه استفاده از هوش مصنوعی و یادگیری

استفاده در سامانه‌های مراکز تماس اطلاع‌رسانی فرودگاهی خواهیم پرداخت؛ به‌گونه‌ای که مخاطبین پس از تماس تلفنی با مرکز تماس سامانه اطلاع‌رسانی فرودگاه‌های کشور و بیان شماره پرواز مورد نظر، از آخرین وضعیت آن پرواز مطلع خواهند شد. بنابراین با توجه به نیاز مبرم به داشتن ارقام درست جهت ارسال به سامانه‌های ثالث به‌منظور دریافت آخرین وضعیت هر پرواز، این فاز به‌عنوان یکی از اصلی‌ترین بخش‌های سامانه هوشمند ارائه اطلاعات پرواز گویا خواهد بود.

جدول ۳: ابزارهای استفاده شده در ParsAirCall

نام ابزار	توضیحات	نسخه
Linux Ubuntu	سیستم عامل	22.04.2 LTS
Python	زبان برنامه‌نویسی	3.10.12
Pandas	کار با انواع داده‌ها	1.5.3
Scikit-learn	پیش‌پردازش داده‌ها و یادگیری ماشین	1.2.2
Matplotlib	بصری‌سازی	3.7.1
Sox	کار با داده صوتی	1.4.1
ffmpeg	کار با داده صوتی	0.2.0
Wandb	نمایش خروجی و لاگ‌ها	0.16.2
Librosa	کار با داده‌های صوتی	0.9.2
PyTorch	یادگیری عمیق	2.1.0
PyTorch Lightning	انعطاف بالاتر کار با پایتورچ	2.1.3
CUDA	ابزار جهت کار با GPU	12.1

۵. ارزیابی روش پیشنهادی و مقایسه با سایر مدل‌ها

نرخ خطای کلمه (WER) یک معیار رایج برای ارزیابی کارایی سیستم‌های شناسایی گفتار خودکار است. این معیار تعداد جایگذاری‌ها، اضافه‌ها، و حذف‌ها را می‌سنجد تا خروجی شناسایی شده را با متن مرجع هماهنگ کند.

$$WER = \frac{S + D + I}{N} \quad (1)$$

به این صورت که در مقایسه متن خروجی با داده‌های واقعی، 'S' تعداد کل جایگزینی‌های انجام شده و 'D' تعداد حذف‌های انجام شده و 'I' تعداد درج‌های انجام شده را نشان می‌دهند. تعداد تجمعی اصطلاحات در داده‌های واقعی نیز با 'N' مشخص می‌شود. از آنجاکه WER مستقل از زبان است، از آن برای ارزیابی کارایی ASR در هر زبانی، از جمله فارسی، می‌توان استفاده کرد.

جدول ۴: ارزیابی مدل‌های نهایی

15. Hsiao, R., et al., *Online Automatic Speech Recognition With Listen, Attend and Spell Model*. IEEE Signal Processing Letters, 2020. **27**: p. 1889-1893.
16. Collobert, R., C. Puhersch, and G. Synnaeve, *Wav2letter: an end-to-end convnet-based speech recognition system*. arXiv preprint arXiv:1609.03193, 2016.
17. Baevski, A., et al., *wav2vec 2.0: A framework for self-supervised learning of speech representations*. arXiv preprint arXiv:2006.11477, 2020.
18. Peng, C.-J., et al., *Attention-based multi-task learning for speech-enhancement and speaker-identification in multi-speaker dialogue scenario*. arXiv preprint arXiv:2101.02550, 2021.
19. Radford, A., et al. *Robust speech recognition via large-scale weak supervision*. in *International Conference on Machine Learning*. 2023. PMLR.
20. Ardila, R., et al., *Common voice: A massively-multilingual speech corpus*. arXiv preprint arXiv:1912.06670, 2019.
21. Kriman, S., et al. *Quartznet: Deep automatic speech recognition with 1d time-channel separable convolutions*. in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2020. IEEE.

ماشین در راهکارهای ارتباطی به‌ویژه مراکز تماس و مراکز تلفن بومی موارد به شرح ذیل پیشنهاد می‌گردد:

- بررسی تأثیر Data Augmentation داده‌های تلفنی در خروجی نهایی
- پیاده‌سازی مدل‌های مبتنی بر یادگیری عمیق بهینه‌تر به‌گونه‌ای که قابلیت اجرا در کلاس‌های مختلف داشته‌باشد.
- جمع‌آوری مجموعه داده‌گان کامل‌تر در فضاهای نویزی و شهری با لهجه‌ها و گویش‌های کامل‌تر

References

1. Deshpande, G., A. Batliner, and B.W. Schuller, *AI-Based human audio processing for COVID-19: A comprehensive overview*. Pattern recognition, 2022. **122**: p. 108289.
2. Agarwal, P., S. Swami, and S.K. Malhotra, *Artificial intelligence adoption in the post COVID-19 new-normal and role of smart technologies in transforming business: a review*. Journal of Science and Technology Policy Management, 2022.
3. Mitreska, M., et al., *Representation Learning for Automatic Speech Recognition: A Review of Speech-to-Text Methods*. 2023.
4. Young, S.J. and S. Young, *The HTK hidden Markov model toolkit: Design and philosophy*. 1993.
5. Veisi, H. and A. Haji Mami, *Persian speech recognition using deep learning*. International Journal of Speech Technology, 2020. **23**: p. 893-905.
6. Shafieian, M., *Hidden Markov model and Persian speech recognition*. International Journal of Nonlinear Analysis and Applications, 2023. **14**(1): p. 3111-3119.
7. Hinton, G., et al., *Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups*. IEEE Signal processing magazine, 2012. **29**(6): p. 82-97.
8. Farahani, M., et al., *Parsbert: Transformer-based model for persian language understanding*. Neural Processing Letters, 2021. **53**: p. 3831-3847.
9. Gonbadi, L. and N. Ranjbar, *Sentiment Analysis of People's opinion about Iranian National Cars with BERT*.
10. Farhadi, A., et al., *Unsupervised Domain Adaptation for image classification based on Deep Neural Networks*. Intelligent Multimedia Processing and Communication Systems (IMPCS), 2023. **4**(1): p. 27-37.
11. Salazar, J., K. Kirchoff, and Z. Huang. *Self-attention networks for connectionist temporal classification in speech recognition*. in *Icassp 2019-2019 ieee international conference on acoustics, speech and signal processing (icassp)*. 2019. IEEE.
12. Perero-Codosero, J.M., et al. *Exploring Open-Source Deep Learning ASR for Speech-to-Text TV program transcription*. in *IberSPEECH*. 2018.
13. Hannun, A., et al., *Deep speech: Scaling up end-to-end speech recognition*. arXiv preprint arXiv:1412.5567, 2014.
14. Amodei, D., et al. *Deep speech 2: End-to-end speech recognition in english and mandarin*. in *International conference on machine learning*. 2016. PMLR.