



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر

گزارش کارآموزی  
محل کارآموزی: شرکت عصر گویش پرداز

نگارش

امیرمحمد بابائی  
شماره دانشجویی: ۹۸۳۱۰۱۱

استاد کارآموزی  
دکتر احمد نیک آبادی

تابستان ۱۴۰۱

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر

گزارش کارآموزی  
محل کارآموزی: شرکت عصر گویش پرداز

نگارش

امیرمحمد بابائی  
شماره دانشجویی: ۹۸۳۱۰۱۱

استاد کارآموزی  
دکتر احمد نیک آبادی

تابستان ۱۴۰۱

# سپاسگزاری

اینجانب امیر محمد بابائی مراتب امتنان و تشکر خود را نسبت به استاد کارآموزی جناب دکتر احمد نیک آبادی که در گذراندن این دوره کارآموزی همواره مرا یاری نموده‌اند، ابراز می‌دارم.

امیر محمد بابائی  
تأیید شده ۱۴۰۱

## چکیده

امروزه، مدل‌های تبدیل گفتار به متن، کاربرد بسیاری در چت‌بات‌ها، تعامل ساده‌تر انسان با رایانه و شبکه‌های پخش آنلاین ویدیو پیدا کرده‌اند. با این حال، مدل‌های فعلی در موقعیت‌های نویزی و دارای کیفیت پایین، عملکرد ضعیفی از خود نشان می‌دهند. یکی از رویکردهای اصلی برای بهبود عملکرد این سیستم‌ها، استفاده از داده‌های کمکی غیر از سیگنال‌های صوتی می‌باشد. از نمونه این نوع داده‌های کمکی، داده‌های تصویری حرکت لب‌های گوینده می‌باشد که توانایی جبران نقص اطلاعات داده‌های صوتی مربوط به گفتار را دارا می‌باشد. رویکرد غالب در این روش، استفاده از روش‌های نظارت‌شده<sup>۱</sup> می‌باشد اما با این حال، به دلیل محدود بودن داده‌های برچسب گذاری شده صوتی-تصویری، روش مناسبی در حال حاضر نمی‌باشند. رویکرد جدید برای دستیابی به نتایج بهتر، استفاده از مدل‌های خود-نظارتی<sup>۲</sup> می‌باشد که توانایی رسیدن به عملکرد مناسب با استفاده از حجم داده برچسب‌گذاری شده کمتر را دارا هستند. مدل ای-وی هیوبرت<sup>۳</sup> نمونه ای از مدل‌های مبتنی بر این رویکرد می‌باشد. علاوه بر این، به دلیل نبود دادگان مناسب برای حل این مساله در فارسی، مراحل جمع‌آوری یک دادگان صوتی-تصویری فارسی با استفاده از ویدیوهای آرشیو سایت تلویزیون نیز در این گزارش، ذکر شده است.

## واژه‌های کلیدی:

بازشناسی گفتار، بازشناسی گفتار به واسطه صوت و تصویر، دادگان صوتی-تصویری

<sup>1</sup>Supervised

<sup>2</sup>Self-Supervised

<sup>3</sup>AV-HuBERT

# فهرست مطالب

صفحه	عنوان
۱	۱ مقدمه
۳	۲ معرفی محل کارآموزی
۴	۱-۲ معرفی شرکت
۴	۲-۲ محصولات شرکت
۵	۳-۲ زمینه‌های فعالیت
۶	۳ فعالیت‌ها و تجربیات کارآموزی
۷	۱-۳ پروژه پرشین ای-وی-اس-آر
۷	۱-۱-۳ مساله بازشناسی گفتار
۸	۲-۱-۳ بررسی پیشینه
۸	۳-۱-۳ مدل خود-نظارتی ای-وی هیوبرت
۱۱	۴-۱-۳ دادگان‌های صوتی-تصویری موجود
۱۲	۵-۱-۳ ایجاد دادگان فارسی
۱۷	۴ نتیجه‌گیری و پیشنهادها
۱۸	۱-۴ نتیجه‌گیری و جمع‌بندی
۱۸	۲-۴ پیشنهادها
۱۹	منابع و مراجع
۲۰	پیوست
۲۱	واژه‌نامه‌ی انگلیسی به فارسی

## فهرست اشکال

شکل	صفحه
۱-۳ معماری مدل ای-وی هیوبرت	۱۰
۲-۳ چهار نمونه از خروجی‌های خط لوله پردازشی مدل سینکنت	۱۴

## فهرست نمادها

مفهوم

نماد



# فصل اول

## مقدمه

صوت یکی از مهم‌ترین حالات انرژی در جهان ما می‌باشد و راه ارتباطی اصلی بسیاری از انسان‌ها و دیگر موجودات، از طریق سیگنال‌های صوتی می‌باشد. به همین دلیل، درک و پردازش این نوع از داده‌ها، اهمیت بسیاری در عصر حاضر برای ما دارا می‌باشد.

علاوه بر این، یکی از بهترین حالات تعامل انسان با رایانه، استفاده از صوت و دستورات گفتاری است. بنابراین، برای دستیابی به چنین قابلیت، نیاز است که گفتار برای رایانه‌ها قابلیت پردازش و درک پیدا کرده و سپس از آن برای برقراری ارتباط راحت‌تر میان انسان و رایانه استفاده کرد.

برای این کار، امروزه سامانه‌ها و مدل‌هایی وجود دارند که صرفاً بر روی قسمت صوتی گفتار متمرکز می‌باشند. این مدل‌ها با اینکه در موقعیت‌های عادی و بدون نویز، به دقت و عملکرد مناسبی دست پیدا کرده‌اند، اما در موقعیت‌های نویزی و با کیفیت پایین، عملکرد نسبتاً ضعیفی از خود نشان می‌دهند و به همین دلیل مدل‌های قابل اتکایی نمی‌باشند.

یکی از راهکارها برای قابل‌اتکا کردن این نوع از مدل‌ها، استفاده از داده‌های کمکی می‌باشد. یک نمونه از این داده‌های کمکی، داده‌های تصویری حرکت لب‌های فرد گوینده می‌باشد. این داده‌ها قابلیت جبران نقص اطلاعات سیگنال‌های صوتی را دارا می‌باشند.

این نوع از عملکرد، معادل عملکرد سیستم شنیداری انسان نیز می‌باشد. در انسان نیز، با اینکه گوش، مهم‌ترین نقش را ایفا می‌کند، اما تنها مولفه نمی‌باشد. این موضوع زمانی واضح‌تر می‌شود که در یک محیط شلوغ، به دنبال درک جملات بیان شده توسط یک گوینده هستیم. در این حالت، حرکت لب‌های فرد در کنار گفتار ضعیفی که از فرد به ما می‌رسد، در کنار هم منجر به درک درست گفتار بیان شده فرد گوینده از سمت ما می‌شود.

علاوه بر این موضوع، برای ساخت و پیاده‌سازی چنین سامانه‌هایی، یکی از مهم‌ترین ارکان، وجود داده‌های آموزشی می‌باشد. این نوع از داده‌ها، در زبان‌هایی نظیر زبان انگلیسی به نسبت، به مقدار بیشتری وجود دارند این در حالی است که در زبان فارسی حجم دادگان‌های موجود به نسبت، کم‌تر می‌باشد. یکی از مواردی که در این گزارش در رابطه با آن صحبت خواهد شد، روش جمع‌آوری و گردآوری یک دادگان صوتی-تصویری برای ارائه و استفاده در حل مساله بازشناسی گفتار به واسطه صوت و تصویر می‌باشد.

در ادامه، در فصل دوم به معرفی شرکت عصرگوش پردازش پرداخته و بخشی از مهم‌ترین محصولات و زمینه‌های فعالیت این شرکت بررسی خواهند شد. در فصل سوم، تجربیات کسب شده در این دوره کارآموزی سه‌ماهه، بیان خواهد شد و برخی از چالش‌ها و راه‌حل‌هایی که در این دوره ارائه شدند، بررسی خواهند شد. در نهایت در فصل چهارم، نتیجه‌گیری مربوط به این دوره کارآموزی بیان خواهد شد و پیشنهادهایی در جهت بهبود مدل و دادگان ارائه شده، ذکر خواهد شد.

## فصل دوم

### معرفی محل کارآموزی

در این قسمت، به طور مختصر، شرکت عصرگوش پردهاز معرفی شده و در ادامه محصولات اصلی شرکت و همچنین زمینه‌های فعالیت این شرکت ذکر خواهند شد.

## ۱-۲ معرفی شرکت

عصر گوش پردهاز (سهامی خاص) فعال‌ترین شرکت در زمینه هوش مصنوعی و پردهاز سیگنال گفتار بوده که فعالیت خود را از ابتدای سال ۱۳۸۲ شروع کرده است. عمده محصولات و خدمات ارائه شده توسط این شرکت برای نخستین بار در کشور و به صورت حرفه‌ای در زمینه‌های پردهاز و تشخیص گفتار بوده است. این شرکت با پشتوانه فنی گروهی از متخصصان کشور از دانشگاه صنعتی شریف تأسیس شد که سابقه و تجربه پژوهشی آنها در زمینه‌های مرتبط با پردهاز سیگنال به چندین سال قبل از شروع رسمی فعالیت شرکت برمی‌گردد.

## ۲-۲ محصولات شرکت

عصرگوش پردهاز پیشرو در ارائه سیستم‌های مبتنی بر گفتار برای زبان فارسی، محصولات مختلفی را توسعه داده است که بیشتر آنها برای نخستین بار برای زبان فارسی انجام شده و منحصرأ توسط این شرکت تولید می‌شوند. برخی از محصولات این شرکت عبارتند از:

- نویسا: نخستین سامانه تایپ گفتاری فارسی
- نیوشا: نخستین سامانه تلفن گویای هوشمند مبتنی بر گفتار
- آریانا: سامانه متن به گفتار فارسی با صدای طبیعی
- شناسا: تعیین هویت گوینده
- رمزآوا: احراز هویت گوینده
- بینا: تصویر خوان هوشمند
- رومند: چت بات هوشمند
- جویا: سامانه جستجوی عبارات و کلمات در گفتار
- پوشا: سامانه پنهان سازی اطلاعات در تصویر (استگانوگرافی)
- پدیدا: سامانه کشف تصاویر پنهان نگاری شده
- پارسیا: اولین نرم‌افزار مترجم گفتار به گفتار فارسی به انگلیسی/عربی

- نویسیار: اولین نرم افزار تایپ هوشمند فارسی
- کارا: نخستین سامانه تشخیص فرمان صوتی برای ویندوز

## ۳-۲ زمینه های فعالیت

این شرکت امروزه دارای گروهی متخصص و منسجم از افرادی با تخصص و تجربه بالا بوده و سابقه طولانی و موفق در زمینه تحقیق و توسعه و کاربردی کردن توانمندی های پژوهشی دارد و علاوه بر ارائه محصولات مختلف در زمینه های هوش مصنوعی، پردازش گفتار فارسی و انگلیسی و پردازش تصویر، قادر به انجام پروژه های مختلف و ارائه خدمات در زمینه های مختلف نرم افزاری می باشد. از جمله زمینه های فعالیت این شرکت:

- تولید نرم افزارها و سخت افزارهای هوشمند
- هوش مصنوعی و شناسایی الگو
- پردازش سیگنال (گفتار و تصویر)
- تشخیص گفتار و تایپ گفتاری (تبدیل گفتار به متن)
- سنتز گفتار و متن خوان (تبدیل متن به گفتار)
- شناسایی افراد از روی صدا
- پردازش زبان طبیعی
- بهبود کیفیت گفتار
- طراحی دادگان های گفتاری و متنی
- طراحی، توسعه و پشتیبانی نرم افزارهای کاربردی مرتبط
- سیستم های تلفن گویا (با قابلیت تشخیص گفتار)
- سامانه های تلفنی مبتنی بر ویپ (استریسک، الستیکس و ...)
- برنامه نویسی روی ریز کامپیوترها (DSP، تلفن همراه و ...)

با توجه به نوآوری های انجام گرفته در شرکت عصرگویش پرداز، این شرکت علاوه بر انتشار مقاله های مختلف در نشریات و کنفرانس های علمی ملی و بین المللی، دارای افتخارات و تأییدیه های متعددی می باشد.

## فصل سوم

### فعالیت‌ها و تجربیات کارآموزی

در این قسمت به تجربیات کسب شده در دوره کارآموزی شرکت عصرگوش پرداز پرداخته خواهد شد. در این دوره کارآموزی، در پروژه بازشناسی گفتار به واسطه صوت و تصویر (پرشین ای-وی-اس-آر<sup>۱</sup>) فعالیت داشته‌ام. در ادامه، فعالیت‌های انجام شده در این پروژه به تفصیل بیان خواهد شد.

### ۱-۳ پروژه پرشین ای-وی-اس-آر

در این بخش، در ابتدا به صورت خلاصه مساله و ضرورت حل آن بررسی خواهد شد سپس به بررسی فعالیت‌های انجام شده در جهت حل این مساله و آماده‌سازی یک خدمت<sup>۲</sup> برای ارائه آن، پرداخته خواهد شد.

#### ۱-۱-۳ مساله بازشناسی گفتار

مهم‌ترین راه ارتباطی انسان، زبان و یکی از ارکان مهم آن، گفتار می‌باشد. بنابراین یکی از مناسب‌ترین روش‌ها برای ارتباط و تعامل با رایانه‌ها، گفتار می‌باشد. به همین دلیل این مساله، یکی از مهم‌ترین مسائل عصر حاضر می‌باشد.

رویکرد غالب در جهت حل این مساله، ایجاد سامانه‌ای است که با دریافت گفتار به صورت سیگنال‌های صوتی، آن را درک کند و سپس متن متناظر با گفتار را به عنوان خروجی، برگرداند. این رویکرد، عملکرد مناسبی در موقعیت‌های بدون نویز از خود نشان می‌دهد اما در صورت قرارگیری در محیط‌ها و موقعیت‌های نویزی، دچار افت کیفیت شده و عملکرد ضعیفی از خود نشان می‌دهند [۱].

برای حل این مساله دو رویکرد عمده وجود دارد:

• تقویت گفتار<sup>۳</sup>

• بازشناسی گفتار با استفاده از ترکیب داده‌های صوتی و بصری<sup>۴</sup>

در این پروژه، برای افزایش پایداری<sup>۵</sup> مدل‌های بازشناسی گفتار در محیط‌های نویزی، از رویکرد دوم استفاده شده است. در این رویکرد، مدل تلاش می‌کند با استفاده از داده‌های بصری - به خصوص حرکت لب‌های فرد گوینده - ضعف قسمت‌های نویزی سیگنال‌های صوتی را جبران کرده و عملکرد بهتری از خود نشان دهد.

<sup>۱</sup>PersianAVSR

<sup>۲</sup>Service

<sup>۳</sup>Speech Enhancement (SE)

<sup>۴</sup>Audio-Visual Speech Recognition (AVSR)

<sup>۵</sup>Robustness

## ۳-۱-۲ بررسی پیشینه

اولین گام در این دوره کارآموزی، مرور سوابق پژوهشی در جهت حل این مساله بوده است. برای یافتن مقالات مربوط به این مساله، با استفاده از سایت پیپرزویدکد<sup>۶</sup>، گوگل اسکولار<sup>۷</sup> و کانکتدپیپرز<sup>۸</sup> فرایند جستجو مقالات را آغاز کرده و در نهایت مقالات مرتبط را با در نظر گرفتن پارامترهای زمان انتشار، وجود پیاده‌سازی در سایت گیت‌هاب<sup>۹</sup> و وجود مدل‌های آماده، جمع‌آوری و در یک برگه گوگل ذخیره کردم. لیست مقالات جمع‌آوری شده، در پیوست قابل مشاهده می‌باشد.

پس از جمع‌آوری تمام مقالات، برای یافتن مقاله مناسب، به بررسی تمام مقالات پرداختم. در کل، یازده مقاله جمع‌آوری شده، دارای پیاده‌سازی با استفاده از چارچوب‌های<sup>۱۰</sup> تنسورفلو<sup>۱۱</sup> و پایتورچ<sup>۱۲</sup> بودند. از میان این یازده مقاله، به دلیل جدیدتر بودن، وجود پیاده‌سازی در گیت‌هاب و وجود مدل‌های آماده، مدل ای-وی هیوبرت<sup>۱۳</sup> و مقاله‌های مربوط به آن ([۱] و [۲]) را انتخاب نمودم.

## ۳-۱-۳ مدل خود-نظارتی ای-وی هیوبرت

مدل ای-وی هیوبرت، یک مدل خود-نظارتی<sup>۱۴</sup> می‌باشد و آموزش آن شامل دو مرحله پیش‌آموزش بر روی داده‌های بدون برچسب و کوک کردن آن با استفاده از داده‌های برچسب‌گذاری شده می‌باشد. به همین دلیل، این مدل با استفاده با حجم کم‌تری از داده‌های برچسب‌گذاری شده، عملکرد بهتری نسبت به مدل‌های نظارت‌شده<sup>۱۵</sup> از خود نشان می‌دهد [۱].

ساختار یادگیری این مدل، از رویکرد اصلی آموزش در مدل زبانی معروف برت<sup>۱۶</sup> الهام گرفته شده است. مدل زبانی برت، یک مدل مبتنی بر ترنسفورمر<sup>۱۷</sup> می‌باشد و برای یادگیری سعی می‌کند قسمتی از جمله ورودی - برای مثال تعدادی از کلمات موجود در جمله - را پوشانده و در ادامه با کمک کلمات مجاور و ساختار جمله، کلمات پوشانده شده را حدس بزند. این روش منجر می‌شود با حجم داده برچسب‌گذاری شده کمتر و داده‌های بدون برچسب، مدل درک مناسبی نسبت به ساختار جملات و جایگاه کلمات در جمله به دست آورد [۳].

<sup>۶</sup>Papers With Code (<https://paperswithcode.com>)

<sup>۷</sup>Google Scholar (<https://scholar.google.com>)

<sup>۸</sup>Connected Papers (<https://connectedpapers.com>)

<sup>۹</sup>Github (<https://github.com>)

<sup>۱۰</sup>Framework

<sup>۱۱</sup>Tensorflow

<sup>۱۲</sup>PyTorch

<sup>۱۳</sup>Audio-Visual HuBERT (AV-HuBERT)

<sup>۱۴</sup>Self-Supervised

<sup>۱۵</sup>Supervised

<sup>۱۶</sup>Bidirectional Encoder Representations from Transformers (BERT)

<sup>۱۷</sup>Transformer



با الهام از این ایده، مدل هیوبرت<sup>۱۸</sup> برای حل مساله بازشناسی گفتار به واسطه صوت<sup>۱۹</sup> پیشنهاد شده است. یکی از تفاوت‌های اساسی حوزه صوت و متن، ساختار داده ورودی می‌باشد. در حوزه متن، ورودی‌ها به دلیل گسسته بودن، قابل شکستن به ساختارهای کوچک‌تر با معنی به صورت توکن<sup>۲۰</sup> یا کلمات می‌باشند در صورتی که صوت، دارای ماهیت پیوسته بوده و به همین دلیل به طور مستقیم چنین امکانی در این حوزه وجود ندارد. برای حل این مشکل و گسسته سازی صوت و گفتار، پژوهشگران از آواها و هجاها به عنوان کوچکترین ساختارهای معنی‌دار در این حوزه استفاده کرده و گفتار را به صورت ترکیبی از آنها تعریف کردند [۳].

در این مدل، برای استخراج آواها، از استخراج‌کننده ویژگی<sup>۲۱</sup> ام-اف-سی-سی<sup>۲۲</sup> استفاده شده است. این استخراج‌کننده ویژگی با دریافت گفتار، ویژگی‌های با ابعاد ۳۹ را در هر لحظه تولید می‌کند. در نهایت با یک الگوریتم خوشه‌بندی نظیر کا-مینز<sup>۲۳</sup> آواهای اصلی مشخص شده و در فرایند آموزش به عنوان واحدهای سازنده گفتار، شرکت می‌کنند. فرایند استخراج ویژگی، تنها در دور<sup>۲۴</sup> اول به واسطه استخراج‌کننده ویژگی ام-اف-سی-سی انجام شده و در مراحل بعدی به واسطه بازنمایی موجود در لایه‌های میانی شبکه کدکننده<sup>۲۵</sup> ترنسفورمر انجام می‌شود [۳].

در ادامه برای یادگیری مدل، بخشی از آواها و هجاهای اصلی که در فرایند خوشه‌بندی مشخص شده‌اند، در گفتار ورودی پوشانده شده و مدل تلاش می‌کند تا با توجه به ارتباط میان آواها و یادگیری ساختار آنها، بخش پوشانده را حدس بزند. در این روش، از تابع خطا آنتروپی متقاطع<sup>۲۶</sup> و الگوریتم‌های بهینه‌سازی نظیر الگوریتم آدام<sup>۲۷</sup> استفاده شده است [۳].

مدل ای-وی هیوبرت، بر پایه مدل هیوبرت ارائه شده است و رویکردی مشابه را این بار برای حل مساله بازشناسی گفتار به واسطه صوت و تصویر در پی می‌گیرد. همانطور که در تصویر ۱-۳ مشاهده می‌شود، در این مدل فریم‌های صوتی و بصری ویدیو، به ترتیب به واسطه کدکننده صوتی و کدکننده بصری به یک بازنمایی متراکم<sup>۲۸</sup> تبدیل می‌شوند [۲].

در کدکننده بصری، از یک مدل رزنت-هجده<sup>۲۹</sup> شده است. مدل‌های رزنت، به جای ساختار ترتیبی لایه‌ها، دارای اتصالات خارج از ترتیب بوده که موجب کاهش مشکل محوشدن گرادیان<sup>۳۰</sup> و به تبع آن، افزایش تعداد لایه‌های مدل می‌شود. این نوع از مدل‌ها، پیش از ارائه مدل‌های مبتنی بر ویژن ترنسفورمر

<sup>18</sup>HuBERT

<sup>19</sup>Automatic Speech Recognition

<sup>20</sup>token

<sup>21</sup>Feature Extractor

<sup>22</sup>Mel-Frequency cepstrum coefficients (MFCC)

<sup>23</sup>K-means

<sup>24</sup>Epoch

<sup>25</sup>Encoder

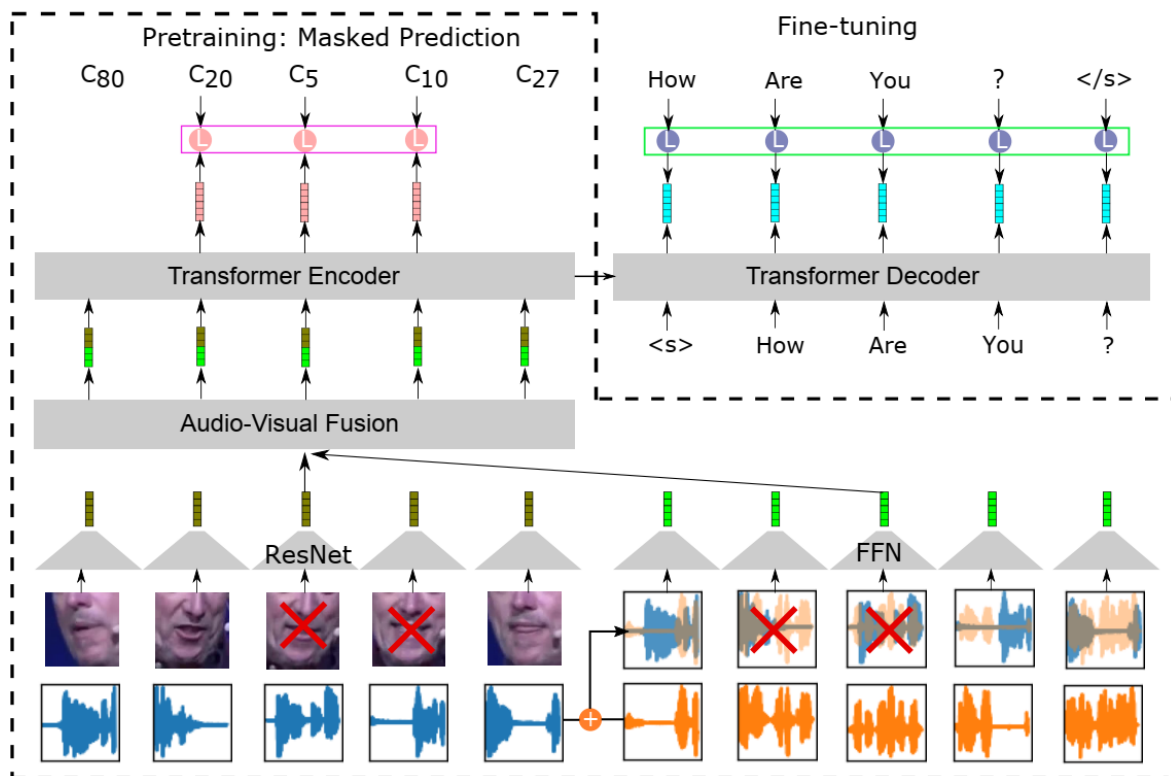
<sup>26</sup>Cross Entropy

<sup>27</sup>Adam

<sup>28</sup>Dense Representation

<sup>29</sup>ResNet-18

<sup>30</sup>Vanishing Gradient



شکل ۱-۳: معماری مدل ای-وی هیوبرت [۲]

<sup>۳۱</sup>، دارای بهترین عملکرد در حوزه تصویر بوده‌اند. پیش از ورودی دادن فریم‌های بصری ویدیو به این مدل رزنت-هجده، تغییرات زیر بر روی تصویر اعمال می‌شود [۱].

- ۶۸ نقطه کلیدی چهره تشخیص داده شده و سپس به واسطه یک تبدیل خطی، این نقاط کلیدی به یک دستگاه مختصات متمرکز بر چهره انتقال پیدا می‌کند.
- یک منطقه مورد علاقه <sup>۳۲</sup> با ابعاد  $96 \times 96$  حول دهان فرد در چهره بریده می‌شود.
- کانال‌های رنگی تصویر به سطح خاکستری منتقل می‌شوند.
- در جهت داده‌افزایی <sup>۳۳</sup> یک کادر با ابعاد  $88 \times 88$  به صورت تصادفی از منطقه مورد علاقه بریده شده و به صورت تصادفی به صورت افقی قرینه <sup>۳۴</sup> می‌شود.

به دلیل تاثیر بیشتر داده‌های صوتی نسبت به داده‌های بصری در این مساله، برای کاهش تاثیر داده‌های صوتی و افزایش تاثیر داده‌های بصری در یادگیری مدل، از یک شبکه تماماً متصل عصبی استفاده شده است. فریم‌های صوتی خام ورودی، پیش از ورودی داده شدن به شبکه عصبی، به واسطه

<sup>۳۱</sup>Vision Transformer (ViT)

<sup>۳۲</sup>Region of Interest (ROI)

<sup>۳۳</sup>Data Augmentation

<sup>۳۴</sup>Horizontal Flip

یک استخراج کننده ویژگی خاص<sup>۳۵</sup> به یک بردار ۲۶ بعدی با فاصله ۱۰ میلی ثانیه تبدیل می‌شود. علاوه بر این، به دلیل تفاوت نرخ برداشت فریم‌های صوتی و بصری - فریم‌های صوتی با فرکانس صد هرتز و فریم‌های بصری با فرکانس ۲۵ هرتز برداشت می‌شود - به ازای هر فریم بصری، چهار فریم صوتی برداشت می‌شود تا هماهنگی زمانی میان دو نوع داده حفظ شود [۱].

پس از ایجاد یک بازنمایی متراکم از فریم‌های صوتی و بصری ویدیو، این بازنمایی به عنوان ورودی به لایه‌های کدکننده ترنسفورمر داده می‌شود. همانطور که در مدل زبانی برت و مدل بازشناسی گفتار هیوبرت توضیح داده شد، در اینجا نیز شبکه کدکننده به دنبال حدس آواهای پوشانده شده است و تلاش می‌کند با این کار، ماهیت و ارتباط میان آواها را به صورت بدون نظارت درک کند. در این مرحله، با استفاده از داده‌های بدون برچسب، می‌توان پیش‌آموزش<sup>۳۶</sup> مدل را انجام داد و در نهایت برای تکمیل فرایند یادگیری مدل و انتقال دانش کسب شده در فرایند پیش‌آموزش به مساله اصلی، مدل با استفاده از داده‌های برچسب‌خورده، کوک می‌شود. برای این انتقال دانش، از یک شبکه کدبرگردان<sup>۳۷</sup> ترنسفورمری به همراه تابع زیان سی-تی-سی<sup>۳۸</sup> استفاده می‌شود [۱].

### ۳-۱-۴ دادگان‌های صوتی-تصویری موجود

پیش از ارائه مدل‌های خود-نظارتی یا نیمه-نظارت‌شده، رویکرد غالب مدل‌ها در جهت حل مساله بازشناسی گفتار به واسطه صوت و تصویر، مبتنی بر مدل‌های نظارت‌شده بوده است. به همین دلیل، عملکرد و دقت این مدل‌ها به حجم دادگان وابسته بوده و در صورت محدود بودن آن، کیفیت و عملکرد پایینی از خود نشان می‌دادند.

با توجه به این توضیحات، بررسی دادگان‌های موجود ضروری به نظر می‌رسد. در ادامه، به صورت مختصر شرحی در رابطه با مشخصات دادگان‌های ال-آر-اس دو و سه<sup>۳۹</sup> و وکس-سلب<sup>۴۰</sup> بیان خواهد شد.

#### دادگان ال-آر-اس دو و سه

دادگان ال-آر-اس دو، یکی از بزرگ‌ترین دادگان‌های موجود برای مساله لب‌خوانی می‌باشد. این دادگان با استفاده از برنامه‌های تلویزیونی - به ویژه اخبار و برنامه‌های گفتگو محور<sup>۴۱</sup> - شبکه انگلیسی زبان بی-بی-سی<sup>۴۲</sup> تشکیل شده است. پس از ارائه این دادگان، دادگان دیگری با نام ال-آر-اس سه،

<sup>۳۵</sup>Log Filterbank Energy

<sup>۳۶</sup>Pre-Train

<sup>۳۷</sup>Decoder

<sup>۳۸</sup>Connectionist Temporal Classification Loss (CTC Loss)

<sup>۳۹</sup>LRS2-BBC, LRS3-TED

<sup>۴۰</sup>Vox-Celeb

<sup>۴۱</sup>Talk Show

<sup>۴۲</sup>BBC

این بار با استفاده از ویدیوهای سخنرانی‌های برنامه‌های تد<sup>۴۳</sup> و تد-ایکس<sup>۴۴</sup> توسط محققان دانشگاه آکسفورد<sup>۴۵</sup> منتشر شد. این دو دادگان، از بزرگ‌ترین دادگان‌های برچسب‌گذاری شده به زبان انگلیسی می‌باشند و معیار ارزیابی<sup>۴۶</sup> در مساله بازشناسی گفتار به واسطه صوت و تصویر محسوب می‌شوند [۱].

#### دادگان وکس-سلب

این دادگان، یک دادگان چندزبانه است که در ابتدا برای مساله تشخیص گوینده چندزبانه با استفاده از داده‌های صوتی و بصری ارائه شده است. بر روی هم، این دادگان شامل بیش از دو هزار و ۴۴۲ ساعت گفتار از بیش از شش هزار گوینده که از سایت اشتراک‌گذاری ویدیو یوتیوب<sup>۴۷</sup> استخراج شده است، می‌باشد. همچنین، این دادگان شامل زیرنویس و متن اصلی که در ویدیوها بیان می‌شود، نمی‌باشد [۱]. مزیت این دادگان نسبت به دادگان ال-آر-اس دو و سه، تنوع بیشتر موقعیت‌ها و صحنه‌هایی است که وجود دارد.

### ۳-۱-۵ ایجاد دادگان فارسی

با توجه به توضیحات داده شده در بخش قبل، دادگان‌های موجود برای حل مساله بازشناسی گفتار به واسطه صوت و تصویر، غالباً به زبان انگلیسی بوده و در زبان فارسی قابل استفاده نمی‌باشند. به همین دلیل، در این دوره تصمیم به ایجاد و جمع‌آوری یک دادگان فارسی گرفتیم. پیش از شروع فرایند جمع‌آوری ویدیوها، مقالات متناظر با دادگان‌های ال-آر-اس سه و وکس-سلب را مطالعه کردم. با توجه به نکات ذکر شده در این مقالات، می‌بایست به سوالات زیر جواب داده می‌شد:

- منبع جمع‌آوری ویدیوها
- چگونگی استخراج ویدیوها
- چگونگی پردازش ویدیوها
- چگونگی فرایند برچسب‌زنی داده‌های استخراج‌شده

#### منبع جمع‌آوری ویدیوها

با توجه به این موضوع که داده‌های دادگان‌های مطرح انگلیسی، با استفاده از برنامه‌های تلویزیونی و ویدیوهای اشتراک گذاشته شده در سایت اشتراک‌گذاری ویدیو یوتیوب به دست آمده بودند، گزینه‌های زیر از گزینه‌های مطرح برای جمع‌آوری ویدیوهای فارسی بودند:

<sup>43</sup>TED

<sup>44</sup>TEDx

<sup>45</sup>Oxford University

<sup>46</sup>Benchmark

<sup>47</sup>YouTube

- سایت تلویزیون
- سایت آپارات
- سایت یوتیوب
- شبکه اجتماعی توییتر<sup>۴۸</sup>
- شبکه اجتماعی اینستاگرام<sup>۴۹</sup>
- سایت‌های آموزش ویدیویی نظیر فرادرس و مکتب‌خونه

در نهایت، از میان این گزینه‌ها، سایت تلویزیون به دلیل برخورداری از ویدیوهای برنامه‌های تلویزیونی و نیازمندی به پالایش کمتر داده‌های این سایت، انتخاب شد. بنا به الگویی از دادگان ال-آر-اس دو و سه، در این مرحله تصمیم نهایی بر استخراج ویدیوهای آرشیوی شبکه خبر صدا و سیما جمهوری اسلامی ایران شد.

در ادامه برای ارتقای این دادگان، می‌توان از داده‌های ویدیویی دیگر برنامه‌های تلویزیون به خصوص برنامه‌های گفتگو محور دیگر شبکه‌های صدا و سیما جمهوری اسلامی ایران نیز استفاده نمود، اما در این مرحله به دلیل محدود بودن زمان این عمل به نسخه‌های بعدی این دادگان موقوف شده است.

#### چگونگی استخراج ویدیوها

برای استخراج ویدیوهای آرشیو تلویزیون از سایت تلویزیون، یک اسکریپت به زبان پایتون<sup>۵۰</sup> و با استفاده از کتابخانه سلنیوم<sup>۵۱</sup> توسعه داده شده است. زبان پایتون یکی از زبان‌های مفسری<sup>۵۲</sup> مطرح می‌باشد. علاوه بر این، کتابخانه سلنیوم، اجازه استفاده به صورت خودکار از مرورگر را به توسعه‌دهندگان داده و توانایی خودکار سازی فرایندهای مرورگر را با استفاده از زبان‌های برنامه‌نویسی دیگر می‌دهد. از دیگر گزینه‌های مطرح برای انجام فرایند استخراج ویدیوها، استفاده از کتابخانه ریکوئستر<sup>۵۳</sup> و بیوتیفول سوپ<sup>۵۴</sup> بوده است. با این حال، به دلیل بارگذاری کند<sup>۵۵</sup> سایت تلویزیون، امکان استفاده از این کتابخانه‌ها وجود نداشت و به همین دلیل از کتابخانه سلنیوم برای انجام این کار استفاده شده است. این اسکریپت با استفاده از یک برنامه راه‌اندازی<sup>۵۶</sup> مربوط به مرورگر، با مرورگر متصل شده و سپس فرایند اتوماسیون و خودکار سازی را شروع می‌کند. این اسکریپت در ابتدا، با توجه به ورودی تعیین شده

<sup>48</sup>Twitter

<sup>49</sup>Instagram

<sup>50</sup>Python

<sup>51</sup>Selenium

<sup>52</sup>Interpreter

<sup>53</sup>requests

<sup>54</sup>BeautifulSoup

<sup>55</sup>Lazy Loading

<sup>56</sup>Driver

به تعداد روزهای تعیین شده، از آرشیو ویدیوهای امروز شروع کرده و به تدریج به سراغ ویدیوهای روزهای قبل رفته و لینک‌های دانلود مربوط به هر ویدیو را استخراج کرده و در یک فایل متنی، ذخیره می‌کند. علاوه بر این، پس از استخراج لینک‌های دانلود ویدیوها، توانایی شروع دانلود ویدیوها را در اختیار دارد.

#### چگونگی پردازش ویدیوها

ویدیوها پس از دانلود، می‌بایست پردازش شده تا آماده ورودی دادن به مدل شوند. در این قسمت با بررسی اسکریپت‌های موجود به صورت عمومی، متوجه شدیم که خط لوله <sup>۵۷</sup> پردازشی مربوط به دادگان وکس-سلب به صورت عمومی در سایت گیت‌هاب منتشر شده است و بدون تغییر قابل استفاده می‌باشد. این خط لوله، مربوط به مدل صوتی-تصویری سینکنت <sup>۵۸</sup> می‌باشد که برای هماهنگ کردن حرکت لب‌های فرد گوینده و صوت استفاده می‌شود.

در این مدل، با استفاده از چارچوب پایتورچ و مدل‌های تشخیص چهره آماده با نام اس-سه-اف-دی <sup>۵۹</sup> چهره‌های موجود در تصویر تشخیص داده شده و در طول ویدیو دنبال می‌شوند. و در نهایت، تمامی چهره‌های موجود در ویدیو، به صورت ویدیوهایی از ویدیو اصلی جدا شده و برای استفاده در حل مسائل دیگر نظیر مساله بازشناسی گفتار به واسطه صوت و تصویر قابل استفاده می‌باشد. در تصویر ۲-۳ چهار نمونه از خروجی‌های خط لوله پردازشی مدل سینکنت را مشاهده می‌کنید.



شکل ۲-۳: چهار نمونه از خروجی‌های خط لوله پردازشی مدل سینکنت

پیش از استفاده از این خط لوله پردازشی، به دنبال توسعه یک خط لوله از ابتدا بوده‌ام. برای این کار،

<sup>57</sup>Pipeline

<sup>58</sup>SyncNet

<sup>59</sup>S3FD

ابتدا از مدل ام-تی-سی-ان-ان<sup>۶۰</sup> که یک مدل تشخیص چهره با قابلیت آشکارسازی در یک صحنه<sup>۶۱</sup> می‌باشد، استفاده کردم. این مدل قابلیت تشخیص چندین چهره در یک تصویر را داشته است. علاوه بر تشخیص کادر چهره فرد، قادر به تشخیص نقاط کلیدی چهره از جمله محل چشمان، محل لب و محل بینی فرد نیز بوده است.

با این حال، به دلیل دقت بالاتر و سادگی پیاده‌سازی خط لوله مربوط به مدل سینکنت، استفاده از این مدل آماده را برای پردازش نهایی ویدیوها برگزیدم.

البته در فرایند پردازش ویدیوها، دسته‌ای از خروجی‌های این مدل به اشتباه استخراج شده و نیازمند حذف در مرحله بعدی می‌باشند. حالت کلی این ویدیوهای به اشتباه استخراج شده به شکل زیر می‌باشد:

- ویدیوهایی که فرد صحبت نمی‌کند ولی به دلیل وجود گفتار پس‌زمینه به اشتباه استخراج شده است.

- ویدیوهایی که فرد گوینده، ماسک به صورت داشته است.

- ویدیوهایی که دارای صدای دوبله‌شده می‌باشد. برای نمونه ترجمه مترجم بر روی صحبت‌های یک فرد غیر فارسی زبان.

از آنجایی که این نوع از ویدیوهای خروجی، از احتمال رخداد پایینی برخوردارند، می‌توان یکی از دو رویکرد زیر را در مواجهه با آنها در پیش گرفت:

- تشخیص به صورت دستی و حذف آنها

- کوک کردن یک مدل به واسطه تمامی خروجی‌های مدل و تشخیص ویدیوهای با ضریب اعتماد پایین<sup>۶۲</sup>

با توجه به مشورت‌های انجام شده با منتورهای این دوره، تصمیم نهایی بر این شد که رویکرد دوم اتخاذ شود. به گونه‌ای که در ابتدا تمام خروجی‌ها استخراج شده و پس از کوک کردن مدل نهایی، با یافتن ویدیوهای با ضریب اعتماد پایین، این ویدیوها بررسی شده و در صورتی که یکی از حالات ذکر شده در بالا باشند، حذف شوند.

علاوه بر این، به دلیل حجم بالای ویدیوهای دانلود شده (حدود ۲۳۲ گیگابایت)، فرایند پردازش ویدیو به سادگی میسر نبوده است. با توجه به امکانات شرکت، سرور<sup>۶۳</sup> دارای کارت گرافیک<sup>۶۴</sup> شرکت، دارای حافظه کافی برای انتقال کامل داده‌ها به این سرور و سپس پردازش تمامی آن‌ها وجود نداشت. به همین دلیل برای رفع این مشکل، از ارتباط میان سرورهای شبکه به صورت محلی استفاده کرده و در هر مرحله، یک ویدیو را با استفاده از دستور اس-سی-پی<sup>۶۵</sup> از سرور دارای ویدیوها دانلودشده به سرور

<sup>۶۰</sup>MTCNN

<sup>۶۱</sup>Single Shot Detector (SSD)

<sup>۶۲</sup>Low Confident

<sup>۶۳</sup>Server

<sup>۶۴</sup>GPU

<sup>۶۵</sup>scp

فعلی منتقل کرده و پس از اعمال پردازش مربوطه بر روی این ویدیو، خروجی‌های میانی مربوط به این داده را در سرور پردازشی پاک کرده و خروجی نهایی را به سرور اولیه که دارای حجم ذخیره‌سازی کافی بوده است، ارسال می‌کردم.

از آنجایی که فرایند پردازش ویدیوهای دانلود شده فرایندی زمانبر می‌باشد و همچنین ارتباط شبکه محلی میان سرورهای شرکت، دارای سرعت بالایی - حدود ۱ گیگابیت به ازای هر ثانیه - دارا می‌باشد، این عامل منجر به کندی تاثیرگذاری در سیستم نهایی نشده و قابل چشم‌پوشی می‌باشد.

#### چگونگی فرایند برچسب‌زنی داده‌های استخراج‌شده

داده‌های آموزشی مرتبط با مساله بازشناسی گفتار به واسطه صوت و تصویر، علاوه بر فریم‌های صوتی و بصری، نیازمند متن بیان شده در فریم‌های ویدیو نیز می‌باشند. به همین دلیل، نیاز است که پس از دانلود ویدیوها و پردازش آنها، متن مرتبط با هر یک از این ویدیوها استخراج شده و به عنوان برچسب این ویدیو مشخص شود. برای انجام فرایند برچسب‌زنی داده‌ها در دسترس، با توجه به مشاوره‌های انجام شده با منتورهای دوره کارآموزی، راهکارهای زیر مطرح شده است:

- جمع‌سپاری<sup>۶۶</sup> داده‌های ویدیویی
  - برچسب‌گذاری به صورت دستی توسط عوامل دوره کارآموزی
  - استفاده از یک مدل آماده تبدیل گفتار به متن<sup>۶۷</sup> و بررسی خروجی آن به ازای هر یک از ویدیوها
- با توجه به زمان محدود دوره کارآموزی، گزینه اول و دوم به دلیل زمانبر بودن، کنار گذاشته شده و تصمیم نهایی بر این شد که با استفاده از یک مدل تبدیل گفتار به متن و سپس بررسی خروجی این مدل به ازای هر ویدیو، این فرایند را سرعت بخشیده و دادگان را سریع‌تر آماده کرد.
- در حال حاضر، حدود دو هزار و سیصد و نود ویدیو دانلود شده است و این تعداد ویدیو، در حال پردازش توسط خط لوله پردازشی سینک‌نت می‌باشد. به دلیل زمانبر بودن فرایند پردازشی، فعالیت فعلی در این پروژه تا به اینجا محدود شده است. پس از آماده سازی دادگان، امکان کوک کردن مدل ای-وی هیوبرت و دیگر مدل‌های آماده نیز وجود دارد و می‌توان میزان مفید بودن این دادگان فارسی جمع‌آوری شده را به طور بهتری ارزیابی کرد.
- علاوه بر این، این دادگان قابل استفاده در مسائل دیگری نظیر تشخیص لب‌خوانی به زبان فارسی نیز می‌باشد و توانایی استفاده برای کوک کردن مدل‌های آماده موجود را نیز دارا می‌باشد.

<sup>۶۶</sup>Crowd Sourcing

<sup>۶۷</sup>Automatic Speech Recognition



## فصل چهارم

### نتیجه‌گیری و پیشنهادها

در این فصل، در ابتدا به مرور نکات ذکر شده و جمع‌بندی آنها پرداخته و سپس پیشنهادهایی در جهت بهبود و ارتقا سامانه و دادگان جمع‌آوری شده ارائه می‌شود.

#### ۱-۴ نتیجه‌گیری و جمع‌بندی

همانطور که در فصل‌های قبل بررسی شد، استفاده از داده‌های تصویری حرکت لب‌های فرد گوینده، داده مناسبی برای جبران نقص در سیگنال‌های صوتی مربوط به گفتار می‌باشد. علاوه بر این، این مکانیزم در سیستم شنیداری انسان نیز وجود دارد و علاوه بر گوش‌ها، دیدن حرکت لب‌های گوینده نیز تاثیر به سزایی در فهم متن بیان شده توسط گوینده دارد.

همچنین، به دلیل محدود بودن داده‌های برچسب‌گذاری شده ویدیویی، استفاده از رویکردهای خود-نظارتی و نیمه-نظارت‌شده نسبت به رویکرد نظارت‌شده، عملکرد بهتری از خود نشان داده و از پایداری بهتری برخوردار خواهد بود. در این مدل‌ها، تلاش می‌شود که دانش کلی نسبت به ماهیت و ارتباط داده‌ها به دست آمده (در فرایند پیش‌آموزش بر روی داده‌های بدون برچسب) و سپس این دانش به طور خاص بر روی حل مساله مورد نظر کوک شود.

این رویکرد، رویکرد مناسبی برای استفاده در زبان‌هایی است که داده ویدیو کافی نداشته باشند؛ چراکه با وجود داده برچسب‌گذاری شده کم نیز، توانایی رسیدن به عملکرد و دقت مناسب را دارا می‌باشند. با این حال، در زبان فارسی داده ویدیویی مناسب برای این مساله موجود نمی‌باشد. به همین دلیل در این مقاله در پی این بر آمدیم تا دادگان ویدیویی فارسی با استفاده از ویدیوهای آرشیو شبکه خبر صدا و سیما جمهوری اسلامی ایران جمع‌آوری کنیم.

#### ۲-۴ پیشنهادها

از جمله پیشنهادهایی که می‌توان در جهت بهبود دادگان فعلی داد، استفاده از ویدیو برنامه‌های گفتگو محور دیگر شبکه‌های صدا و سیما جمهوری اسلامی ایران می‌باشد. علاوه بر این در صورت توسعه یک خط لوله پردازشی دقیق‌تر برای حذف خروجی‌های اشتباه مدل سینکنت، می‌توان به دادگان با کیفیت بالاتری دست یافت. این دادگان به دلیل حجم تخمینی کمی که دارد، برای فرایند کوک کردن مدل‌های بازشناسی گفتار به کمک صوت و تصویر، مناسب می‌باشد اما برای اجرای فرایند پیش‌آموزش مناسب نمی‌باشد. یکی از دیگر پیشنهادها در جهت بهبود این دادگان، استفاده از دیگر سایت‌های اشتراک‌گذاری ویدیو آنلاین نظیر آپارات، فرادرس و مکتب‌خونه می‌باشد.

## منابع و مراجع

- [1] Shi, Bowen, Hsu, Wei-Ning, Lakhota, Kushal, and Mohamed, Abdelrahman. Learning audio-visual speech representation by masked multimodal cluster prediction. arXiv preprint arXiv:2201.02184, 2022.
- [2] Shi, Bowen, Hsu, Wei-Ning, and Mohamed, Abdelrahman. Robust self-supervised audio-visual speech recognition. arXiv preprint arXiv:2201.01763, 2022.
- [3] Hsu, Wei-Ning, Bolte, Benjamin, Tsai, Yao-Hung Hubert, Lakhota, Kushal, Salakhutdinov, Ruslan, and Mohamed, Abdelrahman. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. CoRR, abs/2106.07447, 2021.

## پیوست

لیست مقالات مرتبط با مساله بازشناسی گفتار به واسطه صوت و تصویر  
[لینک دسترسی به برگه گوگل مقالات جمع‌آوری شده](#)

## واژه‌نامه‌ی انگلیسی به فارسی

A	دور، مرحله Epoch . . . . .
بازشناسی گفتار به واسطه صوت و تصویر Audio-Visual Speech Recognition	کدکننده، رمزکننده Encoder . . . . .
	F
بازشناسی خودکار گفتار Automatic Speech Recognition	کوک کردن Fine-Tune . . . . .
	استخراج‌کننده ویژگی Feature Extractor . . . . .
B	چارچوب Framework . . . . .
Benchmark . . . . . معیار ارزیابی	G
C	کارت گرافیک GPU . . . . .
Cross Entropy . . . . . آنتروپی متقاطع	H
Crowd Sourcing . . . . . جمع‌سپاری	قرینه به صورت افقی Horizontal Flip . . . . .
Confidence . . . . . ضریب اعتماد	I
D	مفسر Interpreter . . . . .
Data Augmentation . . . . . داده‌افزایی	L
Decoder . . . . . کدبرگردان، رمزگشا	بارگذاری کند Lazy Loading . . . . .
Dense . . . . . بازنمایی متراکم، بازنمایی چگال Representation	P
Driver . . . . . برنامه راه‌اندازی	خط لوله Pipeline . . . . .
E	پیش‌آموزش Pre-Train . . . . .
	R

Robustness . . . . . پایداری	Supervised . . . . . نظارت‌شده
Region of Interest . . . . منطقه مورد علاقه	Speech Enhancement . . . . تقویت گفتار
S	T
Self-Supervised . . . . . خود-نظارتی	Talk Show . . . . . برنامه گفتگومحور
Semi-Supervised . . . . . نیمه-نظارت‌شده	Token . . . . . توکن
Service . . . . . خدمت	V
Single Shot . . . . آشکارسازی در یک صحنه	Vanishing Gradient . . . محو شدن گرادیان
Detector	