

چگونگی اعمال pre-trained Wav2Vec2.0 در
تشخیص خود کار گفتار در یک استاندارد ارتباطات
کنترل ترافیک هوایی
دوره کارآموزی عصرگوش پرداز

زهرا رحیمی

۱۴ مرداد ۱۴۰۱



فهرست مطالب

۳	۱	مقدمه
۳	۲	انگیزه انجام پژوهش
۴	۳	دادگان
۵	۴	تشخیص خودکار گفتار
۵	۱.۴	تشخیص خودکار گفتار مبتنی بر روش هیبریدی
۶	۲.۴	تشخیص خودکار گفتار انتها به انتها
۶	۱.۲.۴	تقویت داده (data augmentation)
۶	۲.۲.۴	مدل زبانی (language model)
۷	۳.۲.۴	آموزش افزایشی
۷	۴.۲.۴	ارزیابی streaming
۸	۵	نتایج تجربی
۸	۱.۵	شکستن پارادایم، مبتنی بر روش هیبریدی یا انتها به انتها؟
	۲.۵	آیا داده های افزوده ی جزئی در دامنه (partly-in-domain)، عملکرد سیستم تشخیص گفتار را افزایش می دهد؟
۹		
۱۰	۳.۵	مدل های از پیش آموزش دیده چند زبانه چه کمکی می کند؟
۱۰	۴.۵	برای تنظیم دقیق مدل های Wav2Vec2 و XLS-R به چه مقدار داده نیاز است؟
۱۱	۶	نتیجه گیری

۱ مقدمه

تشخیص خودکار گفتار (ASR) می تواند گفتار را به متون قابل خواندن توسط کامپیوتر ترجمه کند. در کنترل ترافیک هوایی (ATC)، راه اصلی ارتباط بین کنترلرهای ترافیک هوایی (ATCo) و خلبانان گفتار رادیویی است. تشخیص خودکار گفتار برای ترجمه گفتار کنترلرهای ترافیک هوایی و خلبان به سیستم ترافیک هوایی معرفی شده است که می تواند برای کاهش بار کاری و اطمینان از ایمنی پرواز استفاده شود. کارهای اخیر بر روی پیش آموزش تحت نظارت خود (self-supervised pre-training) بر روی داده های گفتاری بدون برچسب در مقیاس بزرگ برای ساخت مدل های صوتی (AM) قوی (E2E) متمرکز شده است که بعداً می توانند در کارهای downstream مانند تشخیص خودکار گفتار تنظیم شوند. ما سناریویی را با تجزیه و تحلیل استحکام مدل های Wav2Vec2.0 و XLS-R در downstream ASR برای حوزه ارتباطات کنترل ترافیک هوایی هدف قرار می دهیم.

ATC با هدایت هواپیما در هوا و روی زمین از طریق ارتباطات صوتی بین کنترلرهای ترافیک هوایی (ATCOs) و خلبانان سر و کار دارد. اینها توسط یک دستور زبان و واژگان کاملاً تعریف شده کنترل می شوند که باید برای فراهم کردن یک پرواز امن و قابل اعتماد از ترافیک هوایی و در عین حال پایین نگه داشتن هزینه های عملیاتی تا حد امکان رعایت می شوند. علی رغم علاقه به یک سیستم تشخیص گفتار خودکار برای ارتباطات کنترل ترافیک هوایی، یک سیستم کاملاً کاربردی در بازار وجود ندارد که از جمله دلایل آن می توان به این دو اشاره کرد:

۱. کارایی های پایین سیستم های معرفی شده (به جای تأخیر در انجام وظایف آنها، بهره وری ها ATCO را افزایش می دهد) در نتیجه اهمیت تشخیص درست دستورات و کاهش خطای تشخیص دستورات خلبان (CA-WER (call-sign word-error-rate))

۲. فقدان داده های گفتاری مشروح در مقیاس بزرگ (کمتر از ۵۰ ساعت داده های گفتاری منبع باز) و هزینه تولید بالای آن، آن را تقریباً غیر عملی می سازد

۲ انگیزه انجام پژوهش

پژوهش های ما سناریوی عدم تطابق دامنه را با پاسخ به سه سوال زیر پوشش می دهد::

عملکرد مدل های E2E از پیش آموزش دیده در حوزه های جدیدی مانند ATC چقدر بی نقص و قدرتمند هستند؟

نتایج ما (جدول ۳) تأیید می کند که مدل های انتها به انتها یا E2E که توسط self-supervised-learning یا SSL پیش آموزش شده اند (pre-trained) (مانند Wav2Vec2) یک نمایش قوی از گفتار را یاد می گیرند. تنظیم دقیق (fine-tune) در یک کار downstream (مانند ASR) از نظر محاسباتی ارزان تر از آموزش از ابتدا است، و برای دستیابی به نتایج قابل مقایسه با ASR مبتنی بر روش هیبریدی، به داده های درون دامنه (in-domain) کمتری نیاز دارد. ما همچنین این فرضیه را مطرح می کنیم که مدل های چند زبانه E2E مانند XLS-R در داده های گفتاری ATC

که حاوی انگلیسی لهجه دار (یعنی مجموعه‌های LiveATC-Test و ATCO2-Test) هستند، بهتر هستند. به دلیل بازنمایی کلی گفتار آموخته‌شده در طول SSL.

چه مقدار داده برچسب‌دار ATC (اعم از صوتی و متنی) در مرحله تنظیم دقیق مورد نیاز است تا نتایج قابل مقایسه با مدل‌های مبتنی بر هیبریدی باشد؟

ما یک مطالعه مقایسه‌ای از ۵ دقیقه (یادگیری چند شات) تا حدود ۱۵ ساعت گفتار برچسب‌دار (یعنی از ۱۰۰ تا ۱۵ هزار گفته) انجام می‌دهیم. علاوه بر این، ما افزایش عملکرد به‌دست‌آمده از رمزگشایی با beam search را با استفاده از یک مدل زبان درون دامنه (LM) به‌جای رمزگشایی بر اساس مدل حریصانه یا greedy بررسی می‌کنیم.

حتی اگر مدل‌های Wav2Vec2 و XLS-R از نظر طراحی قابلیت streaming ندارند، آیا چنین مدل‌های E2E می‌توانند در برنامه‌های بلادرنگ مانند ATC استفاده شوند؟ بسیاری از برنامه‌های کاربردی (به عنوان مثال، ATC) به موتورهای ASR streaming نیاز دارند. ما تاخیر عبور و رمزگشایی هر دو مدل E2E، یعنی Wav2Vec2 و XLS-R را در طول استنتاج ارزیابی می‌کنیم.

۳ دادگان

ما روی دو مجموعه آموزشی و چهار مجموعه تست به زبان انگلیسی با لهجه‌های مختلف آزمایش می‌کنیم (جدول ۱). به این نکته توجه داریم که جمع آوری داده‌های کنترل ترافیک هوایی به دلیل شرایط نويز، حریم خصوصی داده‌ها، میزان گفتار و لهجه زبان، چالش برانگیز و پرهزینه است.

NATS و ISAVIA: داده‌های گفتاری توسط ارائه دهندگان خدمات ناوبری هوایی (ANSP) برای پروژه HAAWAII جمع‌آوری و حاشیه‌نویسی می‌شود. دو مجموعه داده عبارتند از، الف) رویکرد لندن (NATS) و ب) ایسلندی (ISAVIA). در مجموع، ۳۲ ساعت داده‌های رونویسی دستی برای آموزش و ۲ ساعت برای آزمایش وجود دارد. هر دو مجموعه داده به‌عنوان گفتار با کیفیت خوب و با فرکانس ۸ کیلوهرتز فهرست‌بندی می‌شوند. برای دیدن جزئیات بیشتر به جدول ۱ مراجعه کنید.

ATCO2-Test: مجموعه توسعه و ارزیابی موجود به‌عنوان منبع باز و ارائه شده در Interspeech در سال ۲۰۲۱. داده‌ها از ارتباطات کنترل ترافیک هوایی از فرودگاه‌های مختلف واقع در استرالیا، جمهوری چک، اسلواکی و سوئیس تشکیل شده است. این دیتاست حاوی ترکیبی از ضبط‌های پرسر و صدا و با لهجه انگلیسی است. این اولین مطالعه‌ای است که تشخیص خودکار گفتار انتها به انتها (E2E ASR) را برای ATCO2-Test ارزیابی می‌کند، به‌عنوان مثال، نرخ خطای کلمه یا WERهای فهرست‌شده در اینجا می‌توانند به‌عنوان خطوط پایه برای تحقیقات آینده مورد استفاده قرار گیرند.

LiveATC-Test: مجموعه آزمایشی از داده‌های LiveATC4 ضبط‌شده از کانال‌های رادیویی VHF در دسترس عموم، به‌عنوان بخشی از پروژه ATCO2 جمع‌آوری می‌شود و شامل ضبط‌های

آزمایشی و کنترلر ترافیک هوایی با انگلیسی لهجه‌دار از فرودگاه‌های واقع در ایالات متحده، جمهوری چک، است. ایرلند، هلند و سوئیس. ما LiveATC-Test را به عنوان مجموعه داده‌های گفتاری با کیفیت پایین در نظر می‌گیریم، یعنی نسبت سیگنال به نویز (SNR) از ۵ تا ۱۵ دسی‌بل می‌رود. (SNR شاخصی که میزان کیفیت سیگنال را نشان می‌دهد و SNR بالای ۱۰ الی ۱۵ نشان دهنده کیفیت مقبول است)

Characteristics			
Dataset	Train / Test	SNR [dB]	WER [%] [†]
NATS	18h / 0.9h	≥ 20	7.7
ISAVIA	14h / 1h	15-20	12.5
ATCO2-Test	- / 1h	10-15	24.7
LiveATC-Test	- / 1.8h	5-15	35.8

جدول ۱: ویژگی‌های داده‌گان آموزشی و آزمایشی

۴ تشخیص خودکار گفتار

راه اندازی ما بر اساس دو مجموعه داده تنظیم دقیق (fine-tune) شده است. اول، ما از ۳۲ ساعت داده‌های مشروح از NATS و ISAVIA استفاده می‌کنیم و ویژگی‌های آن را در جدول ۱ فهرست می‌کنیم. در حال حاضر از این مجموعه‌های تنظیم دقیق به عنوان مجموعه‌های تنظیم دقیق ۳۲ ساعت و ۱۳۲ ساعت یاد می‌کنیم.

۱.۴ تشخیص خودکار گفتار مبتنی بر روش هیبریدی

همه آزمایش‌ها با جعبه ابزار Kaldi انجام می‌شوند. مدل‌های پایه از شش لایه کانولوشن و ۱۵ شبکه عصبی تاخیر زمانی فاکتوریزه شده (حدود ۳۱ میلیون پارامتر قابل آموزش) تشکیل شده‌اند. ما دستورالعمل استاندارد آموزش زنجیره ای Kaldi's LF-MMI را دنبال می‌کنیم. ویژگی‌های ورودی MFCC با وضوح بالا با میانگین نرمال کپسترال آنالین (CMN) هستند. ویژگی‌ها با i-vectors گسترش یافته‌اند. در هنگام رمزگشایی از مدل زبانی APRAS سه تایی استفاده می‌کنیم. این مدل برای ۵ دوره در ۱۳۲ ساعت گفتار کنترل ترافیک هوایی (شامل NATS و ISAVIA) آموزش داده شده است. WER در آخرین ستون جدول ۱ فهرست شده است.

۲.۴ تشخیص خود کار گفتار انتها به انتها

ما نتایج چهار پیکربندی مدل های Wav2Vec2/XLS-R را گزارش می کنیم که از فرم پلات HuggingFace واکنشی شده اند. از این پس، ما این مدل ها را به عنوان زیر برچسب گذاری می کنیم:

W2v2-B BASE model	<ul style="list-style-type: none">95M parametersPre-trained on train-set 960h LibriSearch
W2v2-L: Large-960h model	<ul style="list-style-type: none">317M parametersPre-trained then fine-tuned with LibriSearch 960h train-set
W2v2-L-60K: Large-960h-LV60K model	<ul style="list-style-type: none">same as w2v2-L but uses LibriSpeech + 60k h from Libri-Light during the pre-training phase
W2v2-XLS-R: XLS-R model	<ul style="list-style-type: none">300M parameters pre-trained on 436k h of publicly available data in 128 languages

جدول ۲: پیکربندی مدل های Wav2Vec2/XLS-R

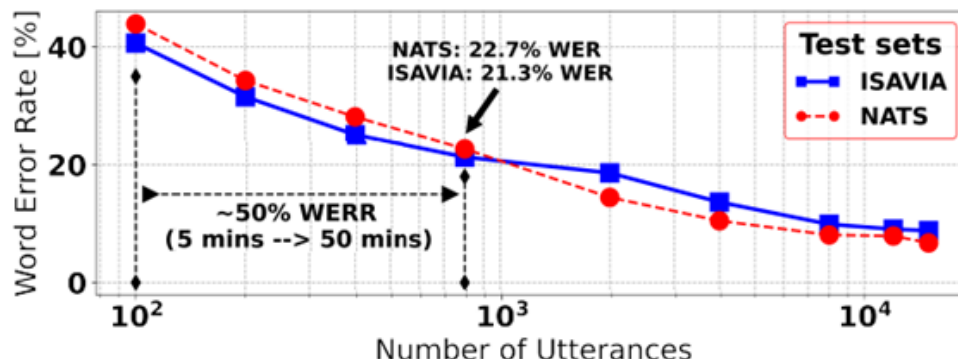
همه آزمایش ها از یک مجموعه فرآپارامتر (hyperparameter) استفاده می کنند. رمز گذار ویژگی (feature encoder) در کل مرحله تنظیم دقیق به روز نمی شود (روش رایج در سناریوهای با منابع کم). ما هر مدل را برای ۱۰ هزار گام با یک مرحله گرم کردن (Warm-up) ۵۰۰ گامی تنظیم می کنیم (تقریباً ۵ درصد از کل به روز رسانی ها). سرعت یادگیری به صورت خطی تا $1e-4$ در طول گرم کردن افزایش می یابد، سپس به صورت خطی تحلیل می رود. ما هر مدل را روی NVIDIA GeForce RTX 3090 با اندازه دسته (batch) ۷۲ (اندازه دسته ۲۴، انباشتگی گرادینان ۳) تنظیم دقیق می کنیم. ما از واژگان مبتنی بر کاراکتر با ابعاد ۳۲ استفاده می کنیم.

۱.۲.۴ تقویت داده (data augmentation)

توالی ورودی را با احتمال $p = 0.075$ و $M = 12$ فریم متوالی ماسک می کنیم. ما همچنین از شاخص فعال سازی ۰.۵۰ استفاده می کنیم. این فرآپارامترها از پیاده سازی اصلی Wav2Vec2 پیروی می کنند.

۲.۲.۴ مدل زبانی (language model)

ما همه رونوشت های متنی را به هم متصل می کنیم و مدله ای زبانی ARPA ۲، ۳ یا ۴ تایی را آموزش می دهیم. مدل های زبانی با همجوشی کم عمق با رمزگشای CTC مبتنی بر پایتون ادغام می شوند. مدل زبانی ۴ تایی به طور سیستماتیک در مقایسه با ۲ تایی ها نتایج بهتری در همه مجموعه های آزمایشی انجام دادند (حدود ۲ درصد WERR نسبی). ما نتایج را فقط با مدل زبانی ۴ تایی گزارش می کنیم. همچنین ما $\alpha = 0.5$ و $\beta = 1.5$ را تنظیم کردیم، که مربوط به طول و وزن نرمال سازی امتیاز مدل زبانی است. اندازه beam را روی ۱۰۰ قرار می دهیم.



شکل ۱: شکل ۱- تاثیر اندازه داده تنظیم دقیق روی نرخ خطای کلمه

۳.۲.۴ آموزش افزایشی

با موفقیت اخیر مدل‌های انتها به انتها که با آموزش تحت نظارت خود از پیش آموزش داده شده‌اند، تعیین مقدار داده‌ای که یک مدل واقعاً برای انجام مؤثر در یک کار downstream نیاز دارد، از اهمیت ویژه‌ای برخوردار است. این امر به ویژه برای کارهای با منابع اندک مانند کنترل ترافیک هوایی که چند ده ساعت داده برچسب گذاری شده برای آموزش یا تنظیم دقیق در دسترس است، بسیار مهم است. در اکثر این نمونه‌ها، داده‌های یک فرودگاه به فرودگاه‌های دیگر به خوبی تعمیم نمی‌یابد و این به دلیل تغییر دامنه قابل توجه AM (لهجه، نرخ بلندگو)، و همچنین تغییر دامنه LM (تسلط هواپیماهای مختلف و دستورات مختلف بسته به فرودگاه) می‌باشد. ما عملکرد مدل را در مقابل اندازه‌های مختلف داده‌های تنظیم دقیق تحلیل می‌کنیم. ما با چهار سناریو یادگیری چند شات با کمتر از یک ساعت (نزدیک ۱۰۰۰ عبارت) داده‌های تنظیم دقیق آزمایش کردیم. در مجموع، ۹ مدل فقط بر روی NATS (خط چین قرمز) یا فقط بر روی داده‌های ISAVIA (خط مستقیم آبی) تنظیم شده‌اند که در شکل ۱ نشان داده شده است (محور x به تعداد عبارت‌های استفاده شده در هنگام تنظیم دقیق در مقیاس لگاریتم اشاره دارد).

هر مجموعه آزمایشی (یعنی NATS و ISAVIA) فقط از داده‌های درون دامنه خود در هنگام تنظیم دقیق و ارزیابی استفاده می‌کند. گفته‌های ۱۰۰، ۱۰۰۰ و ۱۰،۰۰۰ به ترتیب تقریباً ۵ دقیقه (چند شات)، ۱ ساعت و ۱۰ ساعت است.

۴.۲.۴ ارزیابی streaming

مدل‌های Wav2Vec2 و XLS-R با قابلیت‌های streaming طراحی نشده‌اند، اما می‌توان از قابلیت‌های GPU در طول استنتاج برای ارائه رمزگشایی بلادرنگ استفاده کرد. برای آزمایش این فرضیه، روش زیر را انجام می‌دهیم: صوت را به n تکه با اندازه‌های افزایشی تقریباً به اندازه

۳۰۰ میلی ثانیه تقسیم می کنیم. سپس هر تکه افزایشی را به مدل منتقل می کنیم و این کار را برای همه تکه های صوت انجام می دهیم. در نهایت، میانگین زمان مورد نیاز شبکه برای رمزگشایی تمام تکه های یک صوت معین را اندازه گیری می کنیم. ما این فرآیند را روی ۱۰۰ نمونه تصادفی از مجموعه های آزمایشی تکرار می کنیم و میانگین زمان تأخیر را گزارش می کنیم. ما تأثیری را که در راه اندازی streaming بر روی WER ایجاد شده است در نظر نمی گیریم.

۵ نتایج تجربی

در این مقاله، ما فرض می کنیم که مدل های انتها به انتها آموزش دیده شده تحت نظارت خود، یک نمایش قوی از گفتار را یاد می گیرند و در وظایف downstream مانند تشخیص خودکار گفتار تک زبانه یا چند زبانه به خوبی عمل می کنند. ما یافته های خود را با پاسخ به سؤالات زیر تقسیم کردیم:

۱.۵ شکستن پارادایم، مبتنی بر روش هیبریدی یا انتها به انتها؟

اگرچه مدل سازی سیستم تشخیص خودکار مبتنی بر روش هیبریدی برای چندین سال پیش فرض بوده است، اما موج جدیدی از معماری های انتها به انتها که توسط آموزش تحت نظارت خود برای ترکیب اشتراکی AM و LM آموزش داده شده اند جای آن را گرفته اند. ما مدل های انتها به انتها را با بهترین روش مبتنی بر هیبریدی که با مجموعه تنظیم دقیق ۱۳۲ ساعته در Kaldi آموزش داده شده است مقایسه می کنیم (ردیف اول جدول ۳).
برای مدل سازی انتها به انتها:

۱. w2v2-L-60k را برای مجموعه های آزمایشی NATS و ISAVIA انتخاب می کنیم، که فقط در مجموعه ۳۲ ساعته، یعنی داده های درون دامنه، به خوبی تنظیم شد.

۲. سپس w2v2-XLS-R+ برای مجموعه های تست ATCO2-Test و LiveATC-Test، که بر روی ۱۳۲ ساعت از داده های گفتاری کنترل ترافیک هوایی آموزش داده شد، شامل داده های متنوع تر و مشابه مدل مبتنی بر هیبریدی

در نهایت w2v2-L-60k ۳۰ درصد کاهش WER نسبی را در NATS و ۴۱ درصد در ISAVIA در مقایسه با روش مبتنی بر هیبریدی به همراه داشت. این بهبود قابل توجه است، حتی اگر مدل پایه بر روی چهار برابر داده های w2v2-L-60k آموزش داده شده باشد (جدول ۳ را ببینید).
به طور مشابه، w2v2-XLS-R+ (ردیف آخر، جدول ۳) در هر چهار مجموعه آزمایشی از مدل مبتنی بر ترکیبی پیشی می گیرد، اما در ATCO2-Test و LiveATC-Test، دو مورد چالش برانگیز (به دلیل لهجه دار بودن دادگان)، نتایج بسیار قابل توجه است. در مجموع، ۱۹ و ۳۰ درصد WER نسبی در ATCO2-Test و LiveATC-Test به ترتیب به دست آمد (مقایسه w2v2-XLS-R+ نسبت به روش هیبریدی)

Model (num. params.)	Unlabeled data	NATS		ISAVIA		ATCO2-Test		LiveATC-Test		LS*	Latency (ms) [§]
		Greedy	+LM	Greedy	+LM	Greedy	+LM	Greedy	+LM		
Baseline (31M)											
Hybrid-based [†]	-	-	7.7	-	12.5	-	24.7	-	35.8	-	~400 [§]
BASE (95M)											
w2v2-B	LS	10.7	8.4	12.5	10.1	45.6	40.1	48.1	42.2	7.8	32/69
LARGE (371M)											
w2v2-L	LS	9.3	7.6	11.7	9.5	44.9	40.0	47.5	41.4	6.1	33/73
w2v2-L-60k	LS+LV	6.8	5.4	8.8	7.3	34.6	31.2	39.8	34.5	4.9	33/76
w2v2-L-60k+ ^{††}	LS+LV	9.3	<u>7.4</u>	11.2	9.1	23.3	21.2	31.1	27.2	-	-/-
XLS-R (300M)											
w2v2-XLS-R	ML	8.4	6.5	10.5	8.2	39.1	33.8	42.9	36.7	15.4	39/76
w2v2-XLS-R+ ^{††}	ML	<u>9.0</u>	<u>7.4</u>	<u>10.4</u>	<u>8.3</u>	22.8	19.8	29.7	24.9	-	-/-

شکل ۲: WER ثبت شده از چهار مجموعه آزمایشی مطرح شده. WERها به صورت bold و underline به ترتیب مدل‌هایی اشاره می‌کنند که روی داده‌های ۳۲ ساعته و ۱۳۲ ساعته تنظیم شده‌اند.

۲.۵ آیا داده‌های افزوده‌ی جزئی در دامنه (partly-in-domain)، عملکرد سیستم تشخیص گفتار را افزایش می‌دهد؟

ما به این سوال با مقایسه مدل‌های تنظیم‌شده دقیق در مجموعه ۱۳۲ ساعت یا ۳۲ ساعت پاسخ می‌دهیم. توجه داشته باشید که NATS و ISAVIA مجموعه‌های گفتاری کنترل ترافیک هوایی درون دامنه‌ای تمیز هستند، به عنوان مثال، به عنوان درون دامنه برای ۳۲ ساعت و در غیر این صورت تا حدی درون دامنه (in-domain) در نظر گرفته می‌شوند (مجموعه ۱۳۲ ساعت). ATCO2-Test و LiveATC-Test را می‌توان به عنوان مجموعه‌های پر سر و صدا و تا حدی در دامنه (فروگاه‌های مختلف، به عنوان مثال، عدم تطابق صوتی و LM) در نظر گرفت. ما روی w2v2-L-60k و w2v2-L-60k+ که به ترتیب در مجموعه‌های ۳۲ و ۱۳۲ ساعت تنظیم شده‌اند، تمرکز می‌کنیم. توجه داشته باشید که نتایج قابل مقایسه‌ای بین w2v2-XLS-R و w2v2-XLS-R+ وجود دارد. ما WER را در رمزگشایی حریصانه تجزیه و تحلیل می‌کنیم تا فقط روی AM+LM مشترک تمرکز کنیم. در اینصورت یک تنزل WER برای مجموعه‌های آزمایشی درون دامنه، NATS: 6.8% → 9.3% و ISAVIA: 8.8% → 11.2% مشاهده کردیم. این عمدتاً برای افزودن داده‌هایی است که با NATS و ISAVIA مطابقت ندارند. برعکس، کاهش WER قابل توجهی در مجموعه‌های دامنه‌های جزئی وجود دارد، ATCO2-Test: 34.6% → 23.3% و LiveATC-Test 39.8% → 31.1%.

بطور خلاصه مجموعه آزمون NATS با هفت درصد کاهش WER نسبی بدتر نسبت به حالت بدون داده‌های افزوده، تحت تأثیر افزودن داده‌های جزئی در دامنه قرار گرفت (این درصد برای ISAVIA حدود یک درصد است). با این وجود، مجموعه‌های تست چالش برانگیز به طور چشمگیری بهبود یافتند، یعنی ATCO2-Test، ۴۳ درصد و LiveATC-Test، ۳۳ درصد WER نسبی بهبود داشتند.

۳.۵ مدل های از پیش آموزش دیده چند زبانه چه کمکی می کند؟

اگر w2v2-L-60k+ و w2v2-XLS-R را مقایسه کنیم که از تنظیمات دقیق و رمزگشایی beam serach با مدل زبانی استفاده می کنند، نرخ خطای کلمه نسبی در ISAVIA، ATCO2-Test و LiveATC-Test، به ترتیب ۸.۸ درصد، ۶.۶ درصد و ۵.۸ درصد می باشد (بدون پیشرفت در SNR: 5-). در عین حال بهبود قابل توجهی در چالش برانگیزترین مجموعه های تست (10 dB که حاوی گفتار انگلیسی لهجه دار هستند، به عنوان مثال، ATCO2-Test و LiveATC-Test مشاهده می شود. بنابراین، مدل های پیش آموزش شده چند زبانه، در مقایسه با مدل های پیش آموزش شده تک زبانه، عملکرد کمی را افزایش می دهند.

۴.۵ برای تنظیم دقیق مدل های Wav2Vec2 و XLS-R به چه مقدار داده نیاز است؟

ما تأثیر مقادیر مختلف از داده های تنظیم دقیق را در طول مرحله تنظیم دقیق روی نرخ خطاهای کلمه، بررسی می کنیم (شکل ۱). همه آزمایش ها بر اساس قوی ترین مدل انتها به انتها از جدول ۳ هستند، یعنی w2v2-L-60K WER ها با رمزگشایی حریصانه به دست می آیند، یعنی هیچ مدل زبانی یا اطلاعات متنی صریحاً اضافه نمی شود. ما ۱۸ مدل را با تغییر دادن مجموعه داده های آموزشی (یا NATS یا ISAVIA) و مقدار نمونه ها، تنظیم دقیق می کنیم. ما در ابتدا سناریوی یادگیری چند شات (بدترین حالت) را آزمایش کردیم، که در آن تنها ۱۰۰ گفته برجسب دار (۵ دقیقه) برای تنظیم دقیق استفاده شد، و WER ۴۰ درصد برای ISAVIA و ۹.۴۳ درصد برای NATS به دست آورد. علاوه بر این، ۵۰ کاهش WER نسبی با افزایش مقیاس داده های تنظیم دقیق به ۵۰ دقیقه (۸۰۰ گفته) به دست می آید. به طور دقیق، $43.9\% \rightarrow 22.7\%$ NATS و $40.6\% \rightarrow 21.3\%$ WER ISAVIA. در نهایت، اگر از تمام داده های موجود (حدود ۱۴ ساعت) استفاده شود، به ترتیب به ۸.۸ درصد و ۸.۶ درصد WER برای ISAVIA و NATS می رسیم. این نشان دهنده یک WER نسبی ۸۰ درصد در مقایسه با حالت اولیه یا بدترین حالت (۱۰۰ عبارت) است. از دیگر نتایج اینکه با حدود ۸ ساعت (حدود ۸۰۰۰ عبارت) w2v2-L-60K عملکرد تشخیص خود کار گفتار مبتنی بر هیبریدی را شکست می دهد (که از چهار برابر بیشتر داده های آموزشی استفاده می کند).

آیا تشخیص خود کار گفتار real-time در معماری های انتها به انتها، به عنوان مثال، Wav2Vec2 امکان پذیر است؟ ما هر شش مدل از جدول ۳ را در حالت streaming در یک پردازنده گرافیکی NVIDIA GeForce GTX 1080 Ti میان رده آزمایش می کنیم. تأخیرات شامل گذر رو به جلوی مدل، رمزگشایی beam search و رمزگذاری (detokenization) است. نتایج اصلی در جدول ۳ (ستون آخر) گزارش شده است. تأخیر گذر رو به جلو مدل های Wav2Vec2 و XLS-R به طور کلی کمتر از ۱۰۰ میلی ثانیه است. به عنوان مثال، مدل های w2v2-B/L/L- و 60k w2v2-XLS-R زمانی که رمزگشایی حریصانه انجام می دهند تأخیر کمتر از ۴۰ میلی ثانیه دارند. اگر از رمزگشایی beam search با مدل زبانی ۴ تایی استفاده شود، تأخیر تقریباً دو برابر می شود. این تحقیق تخریب WER ناشی از استفاده از مدل های انتها به انتها در حالت streaming

را پوشش نمی دهد.

۶ نتیجه گیری

این مقاله استحکام مدل های Wav2Vec2 از پیش آموزش دیده را در downstream برای کنترل ترافیک هوایی ارزیابی می کند. آزمایش های ما پیشرفت های بزرگی را در تشخیص Wav2Vec2 و XLS-R در مقایسه با تشخیص خود کار گفتار مبتنی بر هیبریدی نشان می دهند. از نظر کمی، بین ۲۰ تا ۴۰ درصد کاهش WER نسبی در مجموعه های آزمایشی ISAVIA، NATS و از مجموعه های چالش برانگیز چند لهجه ای مانند ATCO2-Test و LiveATC-Test به دست آمد. علاوه بر این، ما نشان دادیم که Wav2Vec2 از پیش آموزش دیده، یک مرحله تنظیم دقیق سریع با مقادیر کمی از داده های سازگار را امکان پذیر می کند، به عنوان مثال، تنظیم دقیق ۵ دقیقه ای، مدلی که WER های ۴۰ درصد برای ISVAIA و ۹.۴۳ درصد را برای NATS به دست می آورد. علاوه بر این، ما نشان دادیم که حداقل ۴ ساعت از داده های درون دامنه، WER قابل قبولی در حدود ۱۰ درصد برای ضبط های ISAVIA و NATS ارائه می کنند و با استفاده از داده های دو برابر بیشتر (یعنی ۸ ساعت) عملکردشان از تشخیص خود کار گفتار مبتنی بر هیبریدی پیشی می گیرد. در نهایت، ما اعداد قابل رقابتی در تأخیر برای مدل های Wav2Vec2 و XLS-R در یک GPU میان رده به دست آوردیم، یعنی حدود ۴۰/۸۰ میلی ثانیه با مدل زبانی در رمزگشایی جستجوی حریصانه و beam search.

مراجع

J. Zuluaga-Gomez et al., "How Does Pre-trained Wav2Vec2.0 Perform on [۱] Domain Shifted ASR? An Extensive Benchmark on Air Traffic Control Communications." arXiv, Mar. 31, 2022. doi: 10.48550/arXiv.2203.16822.