



Design and collecting the speaker independent data set of voice commands for controlling the smart home appliances based on Persian speech

---

Leila Safarpour Kalkhoran, Shima Tabibian and  
Elaheh Homayounvala

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

October 12, 2019

# طراحی و جمع‌آوری دادگان مستقل از گویشور حاوی فرامین صوتی برای کنترل لوازم خانگی هوشمند مبتنی بر گفتار فارسی

## برای بیست و پنجمین کنفرانس بین‌المللی انجمن کامپیوتر ایران

لیلا صفرپور کلخوران<sup>۱</sup>، شیما طیبیان<sup>۲</sup> و الهه همایون‌والا<sup>۳</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، پژوهشکده فضای مجازی، دانشگاه شهید بهشتی، تهران

[l.safarpour@mail.sbu.ac.ir](mailto:l.safarpour@mail.sbu.ac.ir)

<sup>۲</sup> استادیار، پژوهشکده فضای مجازی، دانشگاه شهید بهشتی، تهران

[sh\\_tabibian@sbu.ac.ir](mailto:sh_tabibian@sbu.ac.ir)

<sup>۳</sup> استادیار، دانشکده کامپیوتر، گلد اسمیت، دانشگاه لندن، لندن

[e.homayounvala@gold.ac.uk](mailto:e.homayounvala@gold.ac.uk)

### چکیده

امروزه با پیشرفت روزافزون تکنولوژی و سیستم‌های پردازش گفتار، استفاده از گفتار علاوه بر تعاملات انسان‌ها با یکدیگر در تعاملات میان انسان و ماشین نیز ورود پیدا کرده است. یکی از جنبه‌های به‌کارگیری گفتار در زمینه‌هایی مانند خانه‌های هوشمند، جایگزین کردن دستورات صوتی، با فرمان‌هایی است که از طریق لمس به دستگاه داده می‌شود. به این منظور یک سیستم تشخیص فرامین صوتی می‌تواند طراحی و پیاده‌سازی شود. سیستم مذکور قابلیت آن را دارد که هرکجا کاربر درخواستی از سیستم داشت، به راحتی با استفاده از دستورات ساده گفتاری، نیاز خود را برآورده کند. یکی از نیازهای اولیه برای داشتن هر سیستم بازشناس گفتار، وجود دادگان غنی برای آموزش و ارزیابی آن می‌باشد. با توجه به در دسترس نبودن دادگان حاوی فرامین صوتی برای کنترل لوازم خانگی هوشمند به زبان فارسی، در این مقاله به بررسی روند طراحی، جمع‌آوری و ارزیابی یک مجموعه دادگان فرامین صوتی مستقل از گویشور برای کنترل لوازم خانگی هوشمند (تلویزیون، ضبط صوت، لامپ) مبتنی بر گفتار فارسی پرداخته شده است. این دادگان به دو بخش آموزش و آزمون تقسیم شده است. برای ارزیابی دادگان از روش مدل مخفی مارکف و روش حافظه کوتاه مدت ماندگار استفاده شده است. نتایج ارزیابی دادگان مبتنی بر مدل مخفی مارکف نشان می‌دهد که بازشناس کلمات آموزش یافته بر روی دادگان مذکور در بهترین حالت از صحت ۹۳٪ و دقت ۹۲٪ برخوردار بوده و میزان خطا در سطح کلمه تنها ۸٪ است. همچنین نتایج ارزیابی مبتنی بر حافظه کوتاه مدت ماندگار حاکی از آن است که دقت بازشناسی در سطح کلمه با بهترین تنظیمات از میان تنظیمات مختلف انجام شده، حدوداً برابر با ۹۸٪ است.

### کلمات کلیدی

فرامین صوتی مستقل از گویشور، لوازم خانگی هوشمند، دادگان گفتاری فارسی، مدل مخفی مارکف، حافظه کوتاه مدت ماندگار

## ۱- مقدمه

در مقاله [10]، یک سیستم برای کنترل لوازم خانه هوشمند<sup>۱</sup> ارائه شده است که در آن کاربران می‌توانند از دستورات صوتی برای کنترل لوازم خانگی خود استفاده کنند. سیستم پیشنهادی می‌تواند دستورات صوتی را مستقل از ویژگی‌های گوینده مانند لهجه تشخیص دهد. بنابراین یک سیستم بسیار انعطاف پذیر است. فرمان‌هایی که برای ارزیابی این سیستم استفاده شده‌اند، شامل موارد موجود در جدول (۲) می‌باشد که از افراد مختلف در رده‌های سنی مختلف جمع‌آوری شده است.

جدول (۲): دادگان تهیه شده در مقاله [10]

فرمان‌های انواع گروه‌ها				
گروه دسترسی	گروه روشنایی	گروه تهیه هوا	گروه سودمندی	گروه امنیت
باز کن	آشپزخانه	روشن	پرده	کمک
بسته شو	آشپزخانه-خاموش	خاموش	کشیدن-پرده	سکوت
سکوت کن	پذیرایی		آبیاش	
قرار بگیر	پذیرایی-خاموش		خاموش شدن	
	اتاق		آبیاش	
	اتاق-خاموش			

با توجه به مشاهدات و تحقیقات، یکی از چالش‌های اصلی در زمینه فرامین صوتی آموزش سیستم‌های تشخیص فرامین صوتی بر اساس ساختار کلمات و بدون درکی از مفهوم کلمات (متضاد بودن کلمات، هم معنی بودن کلمات و ...) انجام می‌شود. یکی از راهکارهای پیشنهادی نویسندگان این مقاله برای حل این چالش، ترکیب بحث بازنشاسی گفتار با بحث هستان شناسی می‌باشد. در این صورت سیستم تشخیص فرمان صوتی، کلمه را تشخیص داده و کلمه تشخیص داده شده با مفاهیم موجود در پایگاه دانش که به وسیله هستان شناسی ایجاد شده‌اند، مقایسه می‌شود و به دستگاه مربوطه ارسال می‌شود. در واقع با انجام این کار، سیستم مفهوم درخواست کاربر را درک می‌کند و متناسب با آن عمل می‌کند. چالش دیگری که در حوزه جمع‌آوری دادگان می‌توان به آن اشاره کرد، وابسته بودن دادگان به لهجه یا گویش خاص می‌باشد، که برای رفع این چالش، از افراد با گویش‌ها، جنسیت‌ها و رده‌های سنی متنوع استفاده شده است. با توجه به دسترس ناپذیری دادگان حاوی فرامین صوتی مبتنی بر گفتار فارسی در حوزه کنترل لوازم خانگی هوشمند که علاوه بر بحث بازنشاسی گفتار در بحث هستان شناسی هم مورد استفاده قرار بگیرد، بر آن شدیم که به طراحی و جمع‌آوری این دادگان بپردازیم. در این مقاله به معرفی دادگان صوتی مذکور خواهیم پرداخت. ساختار این مقاله به صورت ذیل می‌باشد. در بخش دوم مقاله، روند جمع‌آوری دادگان و شرایط ضبط آن مورد بررسی قرار خواهد گرفت. سپس در بخش سوم، دادگان ارائه شده مورد ارزیابی قرار خواهد گرفت. جمع‌بندی مقاله در بخش چهارم ارائه خواهد شد.

## ۲- روند جمع‌آوری دادگان

در این بخش به طراحی و روند ضبط دادگان حاوی فرمان‌های صوتی مبتنی بر زبان فارسی برای کنترل لوازم خانگی هوشمند (PVC\_SHA)<sup>۲</sup> خواهیم پرداخت.

گفتار یکی از طبیعی‌ترین راه‌های برقراری ارتباط برای انسان‌هاست. ارتباط انسان با ماشین نیز از این امر مستثنی نیست. یکی از حوزه‌های ارتباط افراد با وسایلشان، خانه‌های هوشمند می‌باشد. مقصود از هوشمند شدن خانه، خودکارسازی و تعاملی کردن فعالیت‌هایی است که در خانه انجام می‌شود؛ مانند تهیه مطبوع هوا، سیستم حرارتی، کنترل میزان روشنایی و سیستم امنیتی. تعاملی نمودن این فعالیت‌ها مبتنی بر بازنشاسی گفتار و درک زبان طبیعی توسط لوازم خانه هوشمند می‌باشد [۱]. فرمان‌های صوتی، نقش مهمی را در زندگی امروزه انسان‌ها بازی می‌کنند. سیستم تشخیص فرامین صوتی، در خانه‌های هوشمند، دستیارهای صوتی، کنترل و مدیریت نرم افزارها و جستجوی صوتی مورد استفاده قرار می‌گیرد. در اغلب کاربردهای مذکور، استفاده از گفتار به جای لمس، کلیک کردن یا فشردن دکمه، علاوه بر محبوبیت و سادگی بیشتر، این امکان را فراهم می‌کند که کاربران بدون از دست دادن تمرکزشان، هم‌زمان به کار دیگری هم مشغول باشند [2].

بدیهی است که این قبیل ارتباطات مبتنی بر گفتار نیازمند استفاده از سیستم‌های کارای بازنشاسی و سنتز گفتار می‌باشد. آموزش یک سیستم کارای بازنشاس گفتار در یک کاربرد خاص علاوه بر استفاده از روش‌های نوین پردازش گفتار، نیازمند وجود دادگان غنی متناسب با آن کاربرد می‌باشد.

در حیطه فرمان‌های صوتی، دادگان گفتاری بسیار معدودی در کشور تهیه شده است، که اغلب این دادگان نیز، توسط آزمایشگاه‌های پردازش صوت و گفتار مقیم در دانشگاه‌ها و برای استفاده‌های شخصی تهیه شده‌اند که معمولاً در معرض استفاده عموم قرار نگرفته و درجایی ثبت نشده‌اند. لازم به ذکر است که دادگان موجود در مقالات [3,4]، حاوی فرامین صوتی است که به ترتیب در زمینه ربات اسکات و فضایما، در داخل کشور جمع‌آوری و ارائه شده است. برخلاف کارهای معدودی که در داخل کشور در حیطه فرمان‌های صوتی مخصوصاً در حوزه خانه‌های هوشمند انجام شده‌اند، در خارج از کشور دادگان بسیاری با هدف‌های مختلف بر روی فرامین صوتی مورد استفاده در خانه‌های هوشمند جمع‌آوری شده‌اند [5-8]. از میان کارهای انجام شده می‌توان به مقاله‌های [9,10] اشاره نمود.

در مقاله [9]، به منظور کنترل بهتر لوازم خانگی هوشمند، از فاصله دور و همچنین تشخیص صدای گفتار انسان از صداهای موجود در طبیعت، یک سری فرمان‌های مربوط به کار با لوازم خانگی هوشمند از افراد مختلف در شرایط مختلف با استفاده از میکروفن‌هایی در مسافت‌های مختلف از لوازم خانگی هوشمند، جمع‌آوری شده‌اند. فرامین ضبط شده مطابق جدول (۱) می‌باشد.

جدول (۱): دادگان تهیه شده در مقاله [9]

فرمان‌ها	روشن کن، خاموش کن، روشن شو، رد شو، افزایش دادن، کاهش دادن، بازی کردن، شروع کن، متوقف کن، مکث کن، ادامه دادن، بررسی کردن، خواندن، تغییر دادن، حرکت بعدی، حرکت قبلی
اشیاء	لامپ، لامپ روشنایی، تلویزیون، ضبط صوت، ماشین لباسشویی، ساعت، هوا، تهیه کننده هوا، یخچال، یخچال فریزر، ماشین لباسشویی، ساعت، زمان سنج
موقعیت	آشپزخانه، اتاق خواب، اتاق پذیرایی، اتاق نشیمن، اتاق غذاخوری، حمام

## ۱-۲- شرایط ضبط دادگان

ضبط دادگان (PVC\_SHA) تحت گوشی تلفن همراه با کمک نرم افزاری به نام Voice Recorder (یک نرم افزار مبتنی بر اندروید) صورت گرفته است. برای هماهنگی بین داده های جمع آوری شده، صدای افراد با فرمت Wav، نرخ بیت ۱۶ کیلوهرتز و به صورت تک بانده ضبط شده است. دادگان (PVC\_SHA) شامل فرامین صوتی مطابق جدول (۳) می باشد. همانطور که در جدول مشخص است، این دادگان از ۳۲ کلمه کلیدی و ۶ کلمه غیر کلیدی و ۵۸ فرمان تشکیل شده است.

جدول (۳) : مشخصات مربوط به دادگان

فرمان ها	کلمات کلیدی	کلمات غیر کلیدی
تلویزیون! خاموش شو	یک	است
تلویزیون! قطع شو	دو	ایست
تلویزیون! کافیس	سه	شو
تلویزیون! بسته شو	چهار	کن
تلویزیون! شروع شو	پنج	رو
تلویزیون! باز شو	وصل	بده
تلویزیون! نمایش بده	بخون	
تلویزیون! پخش کن	صفحه	
تلویزیون! روشن شو	صدا	
تلویزیون! صدا رو کم کن	روشن	
تلویزیون! صدا رو معمولی کن	پخش	
تلویزیون! نور صفحه رو کم کن.	نمایش	
تلویزیون! نور صفحه رو زیاد کن	باز	
تلویزیون! نور صفحه رو معمولی کن	شروع	
تلویزیون! نور صفحه رو متوسط کن	بسته	
تلویزیون! صدا رو زیاد کن	متوقف	
تلویزیون! صدا رو متوسط کن	کافی	
تلویزیون! صدا ضعیف است.	قطع	
تلویزیون! صدا خیلی ضعیف است.	خاموش	
تلویزیون! صدا بلند است.	زیاد	
تلویزیون! صدا خیلی بلند است.	کم	
تلویزیون! متوقف شو ،	متوسط	
ضبط صوت! خاموش شو	معمولی	
ضبط صوت! قطع شو	نور	
ضبط صوت! کافی است	سیستم روشنایی	
ضبط صوت! متوقف شو	لامپ	
ضبط صوت! بسته شو	ضبط صوت	
ضبط صوت! شروع شو	تلویزیون	
ضبط صوت! باز شو	بلند	
ضبط صوت! پخش کن	خیلی	
ضبط صوت! روشن شو	تاریک	
ضبط صوت! بخون	ضعیف	
ضبط صوت! صدای کم کن		
ضبط صوت! صدا رو معمولی کن		
ضبط صوت! صدا رو متوسط کن		
ضبط صوت! صدا رو زیاد کن		
ضبط صوت! صدا ضعیف است		
ضبط صوت! صدا خیلی ضعیف است		
ضبط صوت! صدا بلند است		
ضبط صوت! صدا خیلی بلند است		
سیستم روشنایی! روشن شو		
سیستم روشنایی! وصل شو		

سیستم روشنایی! قطع شو  
سیستم روشنایی! خاموش شو  
سیستم روشنایی! نور رو زیاد کن  
سیستم روشنایی! نور رو معمولی کن  
سیستم روشنایی! نور رو متوسط کن  
سیستم روشنایی! نور رو کم کن  
لامپ! روشن شو  
لامپ! وصل شو  
لامپ! قطع شو  
لامپ! خاموش شو  
لامپ! نور رو کم کن  
لامپ! نور رو زیاد کن  
لامپ! نور رو معمولی کن  
لامپ! نور رو متوسط کن  
لامپ! تاریک است  
لامپ! خیلی تاریک است

این فرامین به دو بخش تقسیم شده است: فرامین مستقیم و فرامین غیرمستقیم. فرامین مستقیم، فرامینی هستند که کاربر، درخواست خود را به صورت واضح بیان می کند؛ مانند تلویزیون! خاموش شو. این بخش از فرمان ها در مدت زمان ۱ ساعت و ۴۲ دقیقه و ۳۰ ثانیه از ۳۰ گویشور، شامل ۱۵ نفر مرد و ۱۵ نفر زن با محدوده سنی بین ۱۳ تا ۵۰ سال و با میزان تحصیلات مختلف تهیه شده است. تعداد گویشوران در محدوده سنی ۱۰ تا ۲۰ سال، دو خانم و چهار آقا، در محدوده سنی ۲۱ تا ۳۰، نه خانم و ۱۳ آقا، در محدوده سنی ۳۱ تا ۴۰، یک خانم و در محدوده سنی ۴۱ تا ۵۰، یک خانم بوده است. فرامین غیرمستقیم، فرامینی است که سیستم تشخیص فرمان، به کمک استدلالی که توسط هستان شناسی<sup>۲</sup> انجام می دهد، به نیاز کاربر پاسخ می دهد. مانند لامپ! اتاق تاریک است. این فرامین، در مدت زمان ۱۶ دقیقه و ۲۰ ثانیه از ۲۰ گویشور شامل ۱۰ نفر مرد و ۱۰ نفر زن با محدوده سنی بین ۱۳ تا ۵۵ سال و با میزان تحصیلات مختلف تهیه شده است. تعداد گویشوران در محدوده سنی ۱۰ تا ۲۰ سال، یک خانم، ۲۱ تا ۳۰ سال، هشت آقا و هفت خانم، ۳۱ تا ۴۰، یک خانم و در نهایت در محدوده سنی ۴۱ تا ۵۵، یک آقا و یک خانم بوده است.

همانطور که در جدول (۳) مشاهده می کنید، برای یک عمل مانند خاموش شدن تلویزیون، از فرامین متنوعی استفاده شده است که هدف از این کار، استفاده از دادگان در بحث هستان شناسی برای رفع چالش بیان شده در قسمت مقدمه می باشد.

اطلاعات دقیق گویشوران بر حسب سن، جنسیت، میزان تحصیلات و لهجه در جدول (۴) و (۵) قرار داده شده است. اطلاعات جدول (۴) و (۵) بر اساس جنسیت مرتب شده اند. شماره گویشوران عدد منحصر به فردی بین ۱ تا ۵۰ می باشد.

جدول (۴) : مشخصات مربوط به گویشوران

### فرمان های مستقیم

شماره گویشور	سن	جنسیت	لهجه	میزان تحصیلات
۱	۲۵	مرد	کاشانی	فوق دیپلم
۲	۲۴	مرد	اراکي	لیسانس
۳	۲۸	مرد	مازندرانی	فوق لیسانس
۴	۲۴	مرد	مازندرانی	لیسانس

هرفایل wav به صورت منحصر به فرد نامگذاری می‌شود. این نام مخفف عبارت ضبط شده در هر فایل (به عنوان مثال برای فرمان ( Mediaplayer Mto (turnon) و Sp (مخفف Speaker) و شماره گویشور که عددی بین ۱ تا ۵۰ است، می‌باشد. به عنوان مثال MtoSp50 نام یکی از فایل‌های حاوی دستور (Mediaplayer turn on) می‌باشد.

## ۲-۲- ویرایش و برچسب گذاری

به منظور ویرایش دادگان از نرم افزار CoolEdit استفاده شده است. هر فایل ضبط شده در محیط نرم افزار، باز و نمایش داده می‌شود. چنانچه صداهایی از قبیل کلیک، سرفه، باز و بسته شدن در ... در بخش‌های سکوت فایل وجود داشته باشد، توسط نرم افزار حذف می‌شود. لازم به ذکر است که هنگام ضبط دادگان، از گویشوران خواسته شده بود در محیط تمیز و عاری از هرگونه نویزی، اقدام به ضبط صدا کنند. بنابراین می‌توان ادعا کرد که تمام فایل‌ها بدون هیچگونه نویز قابل اهمیتی جمع‌آوری شده‌اند.

برچسب گذاری دادگان در سطح کلمه و به صورت دستی با استفاده از نرم افزار CoolEdit انجام شده است. برای قسمت‌های سکوت از برچسب sil و برای کلمات غیر کلیدی از برچسب filler استفاده شده است. سایر کلمات برچسب کلمه کلیدی متناظر را خورده‌اند. بنابراین دادگان (PVC\_SHA) حاوی ۲۷۳۷ فایل wav با طول متوسط ۴ ثانیه و ۲۷۳۷ فایل lab (برچسب متناظر با هر فایل) می‌باشد.

## ۳- ارزیابی دادگان

به منظور ارزیابی دادگان (PVC\_SHA) از دو روش مبتنی بر مدل مخفی مارکف و حافظه کوتاه‌مدت ماندگار (LSTM)<sup>۴</sup> استفاده شده است. دادگان به دو بخش آزمون و آموزش تقسیم شده است. مجموعه آموزش حاوی ۱۷۷۷ فایل و مجموعه آزمون حاوی ۹۶۰ فایل می‌باشد، که این فایل‌ها به صورت تصادفی از بین گویشوران انتخاب شده است. بدیهی است که گویشوران در مجموعه آموزش و آزمون مستقل از هم هستند. در ادامه ارزیابی دادگان مبتنی بر هر یک از دو رویکرد مذکور توضیح داده می‌شود.

### ۳-۱- ارزیابی دادگان مبتنی بر مدل مخفی مارکف

برای آموزش بازشناسی فرامین از جعبه ابزار (HTK) استفاده شده است. تعداد مدل‌های مخفی مارکف، ۳۴ مدل به ازای ۳۲ کلمه کلیدی، یک مدل سکوت (sil) و یک مدل کلمه غیرکلیدی (Filler) است. برای هر مدل وضعیت‌های مختلفی (۸، ۱۰، ۱۲، ۱۴، ۱۶، ۱۸) در نظر گرفته شده است. تعداد مخلوط‌های گوسی در هر وضعیت (از ۴ تا ۶۴ مخلوط گوسی) لحاظ شده و حالت بهینه بر روی این دادگان پس از انجام ارزیابی‌های مختلف مشخص شده است. ویژگی‌های استخراج شده از هر فایل، ضرایب مل کپسترام بوده و تعداد مطلوب این ویژگی‌ها پس از انجام ارزیابی‌های مختلف حاصل شده است.

به منظور انتخاب بهینه تعداد ویژگی‌های مل کپسترام، به ازای هر مجموعه مدل مخفی مارکف، دو دسته ویژگی (۱۲ ضریب مل کپسترام و یک ضریب انرژی به همراه مشتقات اول و دومشان (در مجموع ۳۹ ضریب) و ۱۲ ضریب مل کپسترام و یک ضریب انرژی به همراه مشتقات اول و دوم و

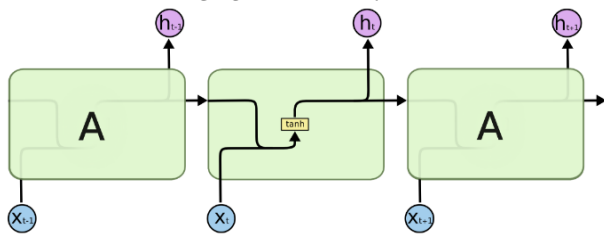
۵	۲۵	مرد	سمنانی	لیسانس
۶	۲۶	مرد	تهرانی	فوق لیسانس
۷	۱۵	مرد	عربی	دانش‌آموز
۸	۲۶	مرد	تهرانی	لیسانس
۹	۲۴	مرد	تهرانی	لیسانس
۱۰	۲۰	مرد	کردی	فوق دیپلم
۱۱	۲۵	مرد	کردی	لیسانس
۱۲	۲۶	مرد	قمی	لیسانس
۱۳	۲۷	مرد	آذری	فوق لیسانس
۱۴	۲۵	مرد	آذری	لیسانس
۱۵	۲۲	مرد	کردی	فوق دیپلم
۱۶	۲۵	زن	آذری	لیسانس
۱۷	۳۱	زن	اراکلی	لیسانس
۱۸	۲۹	زن	شیرازی	فوق لیسانس
۱۹	۲۴	زن	سمنانی	لیسانس
۲۰	۲۴	زن	عربی	لیسانس
۲۱	۳۰	زن	گلستانی	فوق لیسانس
۲۲	۱۹	زن	عربی	دیپلم
۲۳	۲۰	زن	قمی	دیپلم
۲۴	۲۱	زن	سمنانی	دیپلم
۲۵	۲۶	زن	خراسانی	لیسانس
۲۶	۲۶	زن	آذری	لیسانس
۲۷	۱۹	زن	آذری	دیپلم
۲۸	۱۳	زن	آذری	دانش‌آموز
۲۹	۲۲	زن	کردی	فوق دیپلم
۳۰	۵۰	زن	آذری	لیسانس

## جدول (۵) : مشخصات مربوط به گویشوران

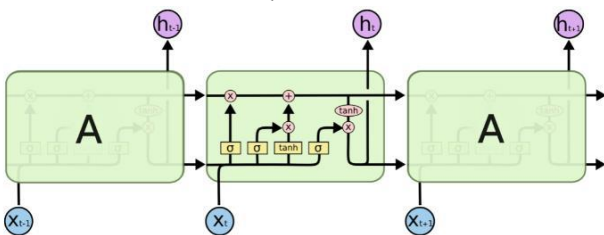
### فرمان‌های غیر مستقیم

شماره گویشور	سن	جنسیت	لهجه	میزان تحصیلات
۱	۲۲	مرد	کردی	فوق دیپلم
۲	۵۵	مرد	آذری	لیسانس
۳	۲۶	مرد	یزدی	لیسانس
۴	۲۴	مرد	آذری	لیسانس
۵	۲۱	مرد	تهرانی	فوق دیپلم
۶	۲۵	مرد	کردی	فوق لیسانس
۷	۳۶	مرد	کرمانی	فوق لیسانس
۸	۲۵	مرد	مازندرانی	لیسانس
۹	۲۸	مرد	مازندرانی	فوق لیسانس
۱۰	۲۷	مرد	گلستانی	لیسانس
۱۱	۲۴	زن	کرمانی	لیسانس
۱۲	۲۸	زن	شیرازی	فوق لیسانس
۱۳	۳۰	زن	یزدی	دکترا
۱۴	۳۱	زن	کردی	دکترا
۱۵	۱۳	زن	آذری	دانش‌آموز
۱۶	۵۰	زن	آذری	لیسانس
۱۷	۲۴	زن	مازندرانی	فوق لیسانس
۱۸	۲۳	زن	خراسانی	لیسانس
۱۹	۲۶	زن	لری	لیسانس
۲۰	۲۵	زن	آذری	لیسانس

حافظه کوتاه مدت ماندگار یا همان LSTM گردید. شکل (۱) تفاوت ساختاری شبکه‌های RNN و LSTM را نشان می‌دهد.



شکل (الف): ساختار RNN



شکل (ب): ساختار LSTM

### شکل (۱): تفاوت ساختار شبکه‌ها [9]

چنانچه شکل (۱) نشان می‌دهد، در مورد شبکه بازگشتی RNN (الف)، ماژول‌های تکرار شونده ساختار بسیار ساده‌ای تنها شامل یک لایه تانژانت هایپربولیک دارند. ورودی هر شبکه، شامل خروجی ماژول قبلی و ورودی جدید می‌باشد. خروجی هر ماژول حاصل اعمال تابع تانژانت هایپربولیک (یا هر تابع فعالیت دیگری) بر ترکیب وزن‌دار دو ورودی مذکور می‌باشد. در شبکه بازگشتی LSTM (ب)، ساختار ماژول تکرار شونده به منظور به خاطر سپاری وابستگی‌های طولانی مدت، پیچیدگی بیشتری دارد. ورودی هر ماژول شامل ورودی جدید و دو خروجی از ماژول قبلی است. همچنین هر ماژول دارای دو خروجی می‌باشد. برای تولید هر یک از این دو خروجی از اعمال توابع فعالیت سیگموئید و تانژانت هایپربولیک (یا هر تابع فعالیت دیگری) بر خروجی‌های ماژول قبلی و ورودی جدید و ترکیب‌های وزن‌دار آن‌ها استفاده شده است. جزئیات مربوط به این بلاک حافظه، اجزای تشکیل دهنده و روابط مربوطه‌اش در [11] ارائه شده اند.

نتایج ارزیابی مربوط به روش LSTM در جدول (۷) نوشته شده است.

جدول (۷): نتایج مربوط به ارزیابی دادگان با LSTM

اندازه کوچک‌ترین بسته <sup>۹</sup>	تعداد واحد در هر لایه مخفی <sup>۷</sup>	بیشترین دوره <sup>۸</sup>	دقت
۲۷	۳۰۰	۳۰۰	%۹۶.۴۱
۳۰	۵۰۰	۵۰۰	%۹۷.۴۵
۳۰	۱۰۰۰	۱۰۰۰	%۹۷.۱۸
۲۰	۱۵۰۰	۱۵۰۰	%۹۷.۲۹
۵	۱۲۰۰	۲۰۰۰	%۹۸.۰۵

با توجه به نتایج حاصل شده در جدول (۷)، بهترین نتیجه، ۲۰۰۰ واحد در هر لایه مخفی و ۱۲۰۰ دوره می‌باشد. این تنظیمات دارای دقت حدود %۹۸، بر بروی دادگان آزمون می‌باشد. همانطور که در بالا اشاره شد، میان دادگان آزمون هیچ گونه اشتراکی از لحاظ گویشوران وجود نداشته و بازشناسی فرامین صوتی، کاملاً مستقل از گویشور عمل می‌کند.

سومشان (در مجموع ۵۲ ضریب)) استخراج شده اند. نتایج ارزیابی در جدول (۶) ارائه شده است.

لازم به ذکر است که ارزیابی‌ها بر اساس سه معیار ارزیابی دقت، صحت و WER (Word Error Rate) انجام شده‌اند که با استفاده از روابط زیر محاسبه می‌شوند [4]:

$$\text{Accuracy} = (H - I) / N \quad (۱)$$

$$\text{Correctness} = H / N \quad (۲)$$

$$\text{WER} = D + S + I / N \quad (۳)$$

که H، معرف تعداد تشخیص‌های درست، N، تعداد کل کلمات قابل شناسایی و I، S و D به ترتیب تعداد خطاهای درج، جایگزینی و حذف کلمه می‌باشند.

جدول (۶): نتایج مربوط به ارزیابی

تعداد مخلوط گوسی	تعداد وضعیت	تعداد ضرایب ویژگی	دقت	صحت	%WER
۱۶	۸	۳۹	%۹۱.۲۳	%۹۲.۵۴	%۱۰.۱۳
۱۶	۱۰	۳۹	%۸۹.۷۵	%۹۱.۲۵	%۱۰.۲۴
۱۶	۱۲	۳۹	%۹۰.۸۸	%۹۲.۰۷	%۹.۱۲
۱۶	۱۴	۳۹	%۹۰.۲۴	%۹۱.۲۳	%۹.۷۵
۱۶	۱۶	۳۹	%۸۹.۸۷	%۹۰.۸۴	%۱۰.۱۲
۱۶	۱۸	۳۹	%۸۹.۰۳	%۹۰.۰۰	%۱۰.۹۶
۴	۸	۳۹	%۸۹.۶۹	%۹۱.۶۶	%۱۰.۳۰
۸	۸	۳۹	%۹۱.۴۱	%۹۲.۸۳	%۸.۵۸
۳۲	۸	۳۹	%۸۷.۷۰	%۸۹.۲۲	%۱۲.۲۹
۶۴	۸	۳۹	%۸۱.۶۶	%۸۳.۴۴	%۱۸.۳۴
۸	۸	۵۲	%۸۱.۲۷	%۸۳.۶۲	%۱۸.۷۳

با توجه به نتایج حاصل شده در جدول (۶)، بهترین دسته مدل و ویژگی، ۸ مخلوط گوسی در هر وضعیت و ۳۹ ویژگی می‌باشد. این تنظیم دارای دقت حدود %۹۲ و صحت %۹۳ و WER حدود ۸٪ بر روی دادگان آزمون می‌باشد. تاکید می‌شود که میان دادگان آزمون هیچ گونه اشتراکی از لحاظ گویشوران وجود نداشته و بازشناسی فرامین صوتی کاملاً مستقل از گویشور عمل می‌کند.

## ۲-۳- مبتنی بر حافظه کوتاه مدت ماندگار (LSTM)

شبکه مبتنی بر حافظه کوتاه مدت ماندگار (LSTM) یکی از انواع شبکه‌های بازگشتی (RNN) می‌باشد که توانایی یادگیری وابستگی‌های بلندمدت را دارد. شبکه LSTM، برای اولین بار توسط هاکریتز و اشمیدر در سال ۱۹۹۷ معرفی شد [10]. در حقیقت هدف از طراحی شبکه LSTM، حل مشکل به یاد سپاری وابستگی‌های بلند مدت در شبکه‌های بازگشتی RNN بود. شبکه‌های بازگشتی RNN ساختاری بسیار مشابه با شبکه‌های چندلایه پرسپترون دارند با این تفاوت که نورون‌های لایه مخفی، علاوه بر یال‌های رو به جلو، یک یال هم به صورت بازگشتی و با احتساب یک زمان تأخیر، از خودشان دارند. چنین ساختاری به یاد سپاری وابستگی‌های کوتاه مدت (Short Term) را تضمین می‌کند. لیکن امکان یادگیری وابستگی‌های مربوط به گذشته‌های دور (Long) را ندارد. برای رفع این مشکل نورون‌های مخفی با یک بلاک حافظه با ساختار پیچیده‌تری، جایگزین شده و منجر به ظهور شبکه‌های مبتنی بر

## ۴- نتیجه گیری

تعامل انسان و دستگاه (HDI<sup>1</sup>) شامل ترجمه دستورات انسان به دستورات کنترلی دستگاه است. علاوه بر این، تعامل یک ارتباط تک جهته نیست، یعنی علاوه بر اینکه دستگاه باید منظور انسان را درک کند، انسان هم باید عملکرد دستگاه را متوجه شود. یک مسئله مهم در HDI این است که مردم به جای کارکردن با رابط‌های سنتی، با استفاده از دستورات خاصی راحت بتوانند اقدامات خود را انجام دهند. برای مثال از طریق یک دستیار صوتی. با اینکه این کار باعث رفاه کاربران می‌شود، اما برای کسانی که با این روش آشنایی ندارند، بسیار آزاردهنده است. موضوعی که از تحقیق‌ها، بررسی‌ها و مشاهدات به دست آمده این است که، افراد برای دستور دادن به دستگاه‌ها از کلمات و عبارات کلیدی خاصی استفاده می‌کنند. به این ترتیب کلمات کلیدی که توسعه دهندگان برای توصیف یک عمل در یک دستگاه استفاده می‌کنند، لزوماً با کلمه‌ای که افراد آشنا به آن وسیله استفاده می‌کنند یکسان نیست و این یکسان نبودن باعث نارضایتی کاربران می‌شود.

برای حل مشکلات مذکور، در این مقاله یک مجموعه دادگان حاوی فرامین صوتی برای کنترل لوازم خانگی هوشمند (لامپ، تلویزیون، ضبط صوت)، ضبط و جمع‌آوری شده است. این دادگان در کار آتی نویسندگان این مقاله در ترکیب هستان شناسی با تشخیص فرامین صوتی، به منظور استفاده از مفهوم کلمه به جای ساختار کلمه با هدف افزایش دقت بازشناسی فرامین صوتی مورد استفاده قرار خواهد گرفت.

چنانچه نتایج ارزیابی دادگان ارائه شده با مدل مخفی مارکف در این مقاله نشان می‌دهند، کارایی کلمه آموزش یافته بر روی دادگان در بهترین حالت از دقت ۹۲ درصد و صحت ۹۳ درصد و ۸ درصد خطای در سطح کلمه برخوردار است.

همچنین نتیجه ارزیابی دادگان ارائه شده مبتنی بر حافظه کوتاه مدت ماندگار با بهترین شبکه حدوداً دارای دقت ۹۸٪ می‌باشد.

## مراجع

[۱] میرحسینی آرانی، سیدمصطفی و مهدی شعبانی، "طراحی رابط صوتی تشخیص فرامین صوتی جهت کنترل سیستم‌های صنعتی"، ششمین کنفرانس ملی علوم و مهندسی کامپیوتر و فناوری اطلاعات، موسسه علمی تحقیقاتی کومه علم آوران دانش، بابل، ۱۳۹۸.

[2] Mehrabani, M., Bangalore, S., & Stern, B.. "Personalized speech recognition for Internet of Things." 2015 IEEE 2nd World Forum on Internet of Things (WF-IoT). IEEE, 2015.

[3] Azargoshasb S., Korayem A. H., Tabibian SH. "A Voice Command Detection system for controlling Movement of SCOUT Robot." 2018 6th RSI International Conference on Robotics and Mechatronics (IcRoM). IEEE, 2018.

[4] Tabibian SH., "A voice command detection system for aerospace applications." International Journal of Speech Technology 20.4 (2017): 1049-1061.

[5] Brenon, A., Portet, F., & Vacher, M. . "Preliminary study of adaptive decision-making system for vocal command in smart home." 2016 12th International Conference on Intelligent Environments (IE). IEEE, 2016.

- [6] Rubio-Drosdov, E., Díaz-Sánchez, D., Almenárez, F., Arias-Cabarcos, P., & Marín, A. "Seamless human-device interaction in the internet of things." IEEE Transactions on Consumer Electronics 63.4 (2017): 490-498.
- [7] Pakpahan, R. L. H., Sudiharto, D. W., & Satwiko, A. G. P. "The prototype of automated doors and windows by using voice commands." 2016 International Seminar on Application for Technology of Information and Communication (ISemantic). IEEE, 2016.
- [8] Mittal, Y., Toshiwal, P., Sharma, S., Singhal, D., Gupta, R., & Mittal, V. K. (2015, December). "A voice-controlled multi-functional Smart Home Automation System." In 2015 Annual IEEE India Conference (INDICON) (pp. 1-6). IEEE.
- [9] Han, Y., Hyun, J., Jeong, T., Yoo, J. H., & Hong, J. W. K. "A smart home control system based on context and human speech." 2016 18th International Conference on Advanced Communication Technology (ICACT). IEEE, 2016. Zachman, John A., "A Framework for Information Systems Architecture", IBM Systems Journal, Vol. 26, No. 3, 1987.
- [10] Arriany, A. A., & Musbah, M. S. "Applying voice recognition technology for smart home networks." 2016 International Conference on Engineering & MIS (ICEMIS). IEEE, 2016.
- [11] Hochreiter, S., & Schmidhuber, J. "Long short-term memory." Neural computation 9.8 (1997): 1735-1780.

## پانویس‌ها

<sup>1</sup> Smart Home Automation System

<sup>2</sup> Persian Voice Commands for Controlling Smart Home Application

<sup>3</sup> Ontology

<sup>4</sup> Long Short Time Memory

<sup>5</sup> Accuracy

<sup>6</sup> Correctnees

<sup>7</sup> Number OF Hidden Unit

<sup>8</sup> MaxEpochs

<sup>9</sup> Mini Bach Size

<sup>10</sup> Human Device intraction