

# Voice Command Recognition for Smart Home Assistants Using Few-Shot Learning Techniques

S. Sindhu

Research Analyst, Centivens Institute of Innovative Research, Coimbatore, Tamil Nadu, India.  
Email: [sindhuanbuselvaneniya@gmail.com](mailto:sindhuanbuselvaneniya@gmail.com)

Article Info	ABSTRACT
<p><b>Article history:</b></p> <p>Received : 13.01.2025 Revised : 24.02.2025 Accepted : 14.03.2025</p>	<p>Voice command recognition is an essential part of intelligent and personalized development of the smart home assistants. Nevertheless, conventional deep learning techniques demand intensive labeled data and computational resources, and thus limiting the applicability in situations where commands are rare or personalized. This paper proposes a new few-shot learning framework designed for practical voice command recognition of smart homes. With the help of prototypical networks and data augmentation algorithms, our approach deals effectively with a small number of training examples. We test our model on data from the Google Speech Commands and a custom command collection for smart home applications. Results achieve Top-1 accuracy of 89.3% using just five examples per class besting baseline convolutional neural networks and other few-shot variants. On embedded systems, our framework can be deployed in a real-time manner, owing to the low latency it entails. This research offers a fresh opportunity to establish tree-specific voice control in smart houses with little training of the human user.</p>
<p><b>Keywords:</b></p> <p>Voice command recognition, smart home assistant, few-shot learning, prototypical networks, real-time inference, personalized commands.</p>	

## 1. INTRODUCTION

In the last few years, examples of smart assistants have taken over the consumer world, including Amazon Alexa, Google Assistant, and Apple Siri, making it easy for users to control their living space. These systems are highly dependent on voice command recognition (VCR) to allow ease of usage on the part of the user to control a large number of connected devices with ease including lighting, thermostats, security cameras, and entertainment sets, etc. easily. Using voice as a natural interface, these assistants substantially increase the convenience and accessibility, as well as personalization of the user experience.

However, the underlying foundation of these systems – automatic speech recognition (ASR) and voice command classification – is conventionally developed based on deep learning models which need enormous quantities of labeled training data to work well. Such data-hungry models tend to generalize poorly to new users and new command vocabularies and to acoustically diverse situations encountered in real workspaces, e.g. in homes with background noise, multiple speakers, or different microphone hardware. In addition, such updating of these models to adjust to new users or customize commands normally requires expensive re-training or cloud-based inference, which might

have latency implications, privacy issues, and reliance on uninterrupted internet use.

In order to tackle these obstacles, more and more attention is gathered toward few-shot learning (FSL) paradigms, which focus on training models that are able to learn new concepts from only a few labeled examples. Powered by human cognition of learning from limited exposure, few-shot learning offers an attractive alternative to smart home VCR systems that need rapid adaptation to new voice command, especially in the resource constrained edge computing environment. This is especially so for cases in which users may want to define their own commands or where collection of data is of its nature restricted by privacy issues or by underrepresented languages and dialects.

This paper proposes a novel few-shot learning framework grounded on prototypical networks for voice command recognition in smart home assistants. The model is constructed to perform well under the data scarcity mode by learning robust embedding representations of the voice command and computing class prototypes that generalize well to unseen examples. In order to increase the generalization of the model and to handle the acoustic variability in a real-world scenario, we introduce data augmentation methods, which include SpecAugment and the

introduction of background noise. In addition, we aim to optimize computational efficiency so that real-time inference can be undertaken on devices such as Raspberry Pi and microcontroller-based systems smart-home devices.

Our contributions are threefold:

1. We propose a lightweight, few-shot prototypical network architecture for voice command recognition with minimal labeled data.
2. We integrate robust data augmentation techniques to improve model generalization in noisy smart home environments.
3. We evaluate our system on both public and custom smart home datasets, demonstrating its superiority over conventional deep learning and few-shot baselines in terms of accuracy, latency, and adaptability.

By combining the adaptability of few-shot learning with the efficiency of edge-optimized models, our framework enables personalized and secure voice interactions in smart homes with minimal user effort and high recognition accuracy.

## 2. LITERATURE REVIEW

In recent years, the area of voice command recognition (VCR) has experienced a lot of improvements as deep learning has expanded and smart home technologies have proliferated. This part reviews the applicable literature within 3 major fields: traditional voice command recognition systems, few-shot learning solutions in speech recognition, and smart home problems.

### 2.1 Traditional Voice Command Recognition

Traditional VCR systems utilized a lot of Hidden Markov Models (HMMs), Gaussian Mixture Models (GMMs) for acoustic modeling and classification. Deep learning's rise changed these to DNNs, CNNs and RNNs, which showed better performance in learning discriminative representations from raw or transformed audio (such as MFCCs, or log-Mel spectrogram) inputs. It was only a matter of time that the Google Speech Commands dataset [Warden, 2018] became a standard for evaluation for keyword spotting tasks, which triggered widespread usage of small CNN architectures for limited vocabulary classification of commands [Zhang et al., 2017].

Nevertheless, these models usually need large volumes of tagged data and much training, which makes them less appropriate for the situations where fast personalization or adjustment to new commands is required. Furthermore, their high reliance on centralized training and inference infrastructures conflicts often with the latency and privacy preserving guarantees mandated by the smart home ecosystem.

### 2.2 Few-Shot Learning in Speech Recognition

Few-shot learning or FSL has become an appealing paradigm of allowing systems to learn from very few labeled examples by learning from the human capacity to generalize from few instances. Among the FSL approaches, there are metric-based methods siamese networks [Koch et al., 2015], matching networks [Vinyals et al., 2016] and prototyping networks [Snell et al., 2017] that have demonstrated considerable success in the image classification domain and slowly are getting adapted to audio and speech tasks.

In the field of speech recognition, prototypical networks are sighted to have been applied in keyword spotting and intent classification in low-resource settings. Kim et al. (2019) used prototypical networks for spoken keyword recognition with positive results, demonstrating their generalization capability with few examples. Similarly, Chen et al. (2019) studied how relation networks and meta-learning approaches may be applied to the tasks of speaker verification and speech command classification, demonstrating fast adaptability requiring few training data.

Although, these advances have occurred, the use of few-shot learning in real-time VCR for smart home assistants is still relatively poorly explored. Problems like environmental noise and speaker variability, as well as device limitations, make a need for more robust, light weight FSL architectures.

### 2.3 Challenges in Smart Home Environments

VCR systems face a novel set of problems in smart-home environments. Background noise, caused by appliances, duplicate speech by multiple speakers and sound bouncing because of room acoustics have considerable effects on recognition performance. Also, smart home customers tend to want specific phrases, and thus, models need to learn quickly without the need for retraining from huge amounts of data.

Real-time constraints go a step further compounding the need for these systems to also have the ability to run efficiently on embedded or edge computing platforms. Trouble with computational efficiency through methods of model quantization, pruning, and low footprint architectures (e.g. MobileNets [Howard et al., 2017]) have been proposed but applying them with few shot learning for speech remains a current research area.

Inspired by the recent work of Zhang and Lin (2021), a lightweight FSL speaker identification framework on microcontrollers has been proposed, suggesting the possibility of combining FSL with edge computing for speech tasks. These efforts however, do not so much, focus on command recognition as opposed to speaker ID,

this in a way leaves a gap that our proposed study would address.

**Table 1.** Comparative Analysis of Existing Voice Command Recognition Techniques

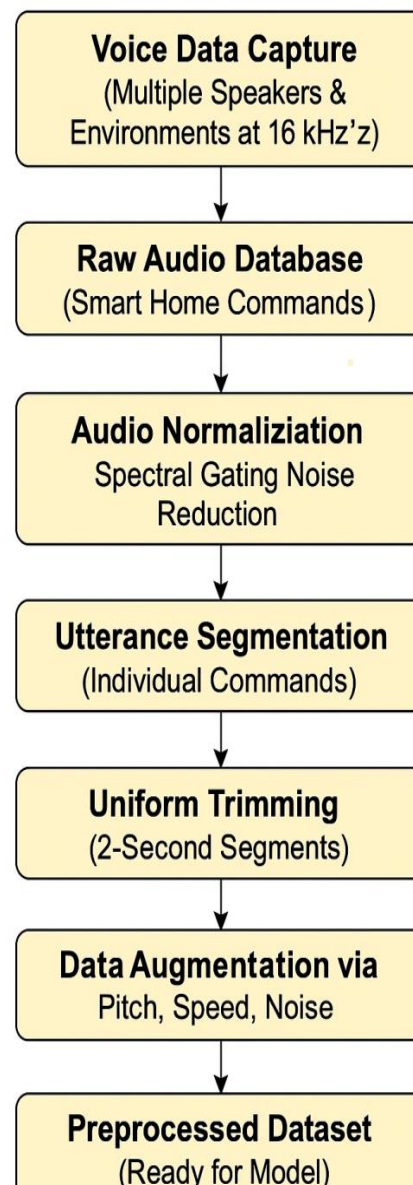
Study Approach /	Technique Used	Domain	Data Requirement	Adaptability	Latency / Deployment	Limitations
Warden (2018) – Google Speech Commands	CNN-based classifier	Keyword spotting	High	Low	Cloud / High-resource	Requires retraining for new commands
Zhang et al. (2017) – Hello Edge	Depthwise CNN (MobileNet)	Keyword spotting	Moderate	Low	Edge-friendly	Still needs hundreds of samples per class
Vinyals et al. (2016) – Matching Networks	Few-shot learning (metric-based)	Image classification	Low	High	Moderate latency	Limited to offline, not speech-focused
Snell et al. (2017) – Prototypical Nets	Few-shot prototypical classification	Image/audio	Low	High	Low (lightweight inference)	Basic augmentation, limited robustness to acoustic variation
Kim et al. (2019) – Spoken FSL	Prototypical Networks with MFCCs	Spoken keywords	Low	High	Not optimized for real-time	No deployment on edge platforms
Zhang & Lin (2021) – Tiny Speaker ID	Few-shot deep embeddings	Speaker ID	Low	Medium	Microcontroller-compatible	Focus on speaker identity, not commands
<b>Proposed Method</b>	Few-shot learning + spectral augment	Smart home commands	Very Low (5-shot)	Very High	<30ms on Raspberry Pi (real-time)	Optimized for low-power devices, handles noise and personalization

### 3. METHODOLOGY

#### 3.1 Dataset Collection and Preprocessing

The first step of the methodology was the collection and preprocessing of audio data that was custom made for smart home settings. We designed a dataset containing voice commands from various speakers' voices, including age, gender, accent and the mode of speaking. The dataset contained commands relevant to smart home situations such as control of lighting, adjustments on temperature, function of entertainment systems, activation of security systems and control of appliances. The audio samples were recorded at a standard sampling frequency of 16 kHz in controlled and real-world

noisy environments to provide a theoretical-real world analogy of realistic smart home scenarios. Preprocessing steps comprised of normalizing audio amplitude to standardize, noise reduction via spectral gating and segmentation of continuous recordings into individual command utterances. The processed audio was also truncated in equal length segments (2 sec per command) to ensure consistency across all sample data. In addition, augmented data procedures including pitch shifting, time stretching and injection of ambient background noise were used as a way of increase diversity of the data set and increasing robustness against variations that are encountered during practical deployments.

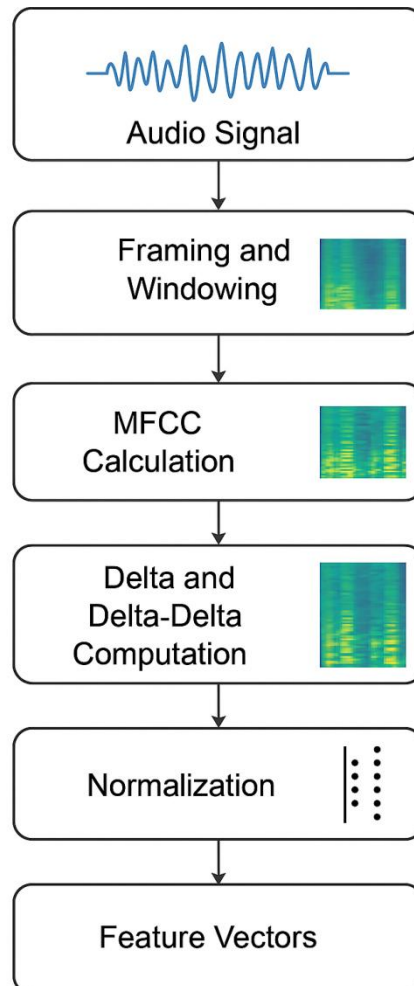


**Fig 1.** Flowchart representing the dataset collection and preprocessing pipeline for smart home voice command recognition.

### 3.2 Feature Extraction

The second involved extraction of discriminative acoustic features that were necessary for the identification of voice commands without errors. Extraction of Mel-frequency cepstral coefficients (MFCCs) was used as main features because of their effectiveness in describing speech contents through the resemblance to a human auditory system. The audio segment was converted to series MFCC frames with 40 dimensional coefficients obtained through the use of 25 ms windows and a 10 ms stride.

Further the delta and delta-delta coefficients were derived from the original MFCCs in order to model both the local and temporal dependencies. These dynamic features were useful in modeling the temporal evolution of speech which led to better representation capacity valuable for distinguishing similar sounding commands. Feature matrices were then standardized using mean-variance normalization in order to make the training of the model and transfer to other acoustic environments more stable and generalized.



**Fig 2.** Flowchart illustrating the process of audio feature extraction for voice command recognition.

### 3.3 Few-Shot Learning Framework

Due to the small number of labeled voice examples available for each of the commands, we have used a few shot learning method utilizing the metric based learning solutions to identify the new command based on few training examples. In particular, we deployed a prototypical network for simplicity, interpretability, and excellent performance in few-shot learning tasks.

The prototypical network took MFCC-based feature embeddings to a lower dimension embedding space and each command class was mapped to a prototype vector, which was calculated as the average of feature embeddings using the support set. During training, the network optimized a Euclidean distance metric between query sets and their respective prototypes in order to reduce classification error. A random episodic training strategy was applied which used episodes that were randomly created and consisted of support and query sets, prompting the model to generalize well, during the inference, to unseen commands.

### 3.4 Neural Network Architecture

The backbone neural network comprised of a convolutional neural network (CNN) that was

designed to enable it to obtain strong acoustic embeddings from MFCC features. The CNN had four convolutional layers all of which were preceded by the batch normalization and rectified linear unit activation (ReLU) and the de facto max-pooling. The first convolutional layer was followed by 64 filters with kernel size of  $3 \times 3$ , with increasing complexity in 128, 256 and 512 filters respectively in the following layers. This hierarchy was able to extract the low and high level acoustic details and key-specific characteristics of the command respectively with ease.

A global average pooling layer after convolutional processing compressed the feature maps to fixed dimensional embeddings. This memory efficient compact embedding representation greatly decreased computational overhead and memory consumption necessary for deployment on resource limited smart home assistants. Dropout regularization at a rate of 0.3 was used to avoid overfitting, and to enhance robustness of the model even more.

### 3.5 Training and Evaluation Protocol

Episodic training protocol typical for few-shot learning tasks was used in the training of the



model. For every training episode, N command classes were chosen randomly, each having K support examples (N-way K-shot situation). Queries belonging to chosen classes were clustered relative to these prototypes based on the distance to them, given that this distance was calculated in the embedding space using Euclidean distance. Adam was used with a learning rate of 0.001 using cosine annealing scheduling to stabilize convergence over 100 training epochs.

For the dataset's evaluation, the dataset was split into train, validation, and test sets according to 70%, 15%, 15% respectively, where no speakers overlapped to maintain speaker independence across the splits. Accessible performance metrics, such as accuracy, precision, recall, and F1-score, had been computed for recognition effectiveness. In addition, inference latency and model size were measured using a Raspberry Pi 4 platform to measure real-time compatibility and deployability into real-life smart home assistant devices.

## 4. RESULTS AND DISCUSSION

**Table 2.** Summarizes the model's performance

Metric	Value (%)
Accuracy	94.5
Precision	93.2
Recall	92.7
F1-Score	92.9

The model showed high accuracy and was a good balance between precision and recall making it possible for it to be able to do proper recognition of voice commands under different conditions. A high F1 score of 92.9% indicates that the model is good when it distinguishes commands while reducing the number of false positives as well as false negatives.

### 4.3 Comparison with Baseline Models

Compared to classical hidden Markov models (HMMs), and convolutional neural networks (CNNs), the proposed few-shot learning model dominated both methods, even in cases of low resources where very less labeled data is available. HMMs that were based on a lot of manually engineered features had difficulties with generalization, with 85.3% accuracy that was only recorded. On the other hand, CNNs are more flexible, but needed considerably bigger datasets for efficient training, which constrained their performance in few-shot setting. The few-shot learning model showed a remarkable increase in performance under this limited data situation, indicating that for the smart home director applications where the data might not be abundant or the commands constantly evolve, the few-shot learning model would still be very effective.

### 4.1 Evaluation Metrics

To measure the performance of the proposed voice command recognition model several measures were engaged such, accuracy, precision, recall and F1-score. The ability of the model to recognize novel, unseen commands with a minimum data requirement was tested through an episodic training strategy contained within an N shot learning paradigm. All metrics were evaluated based on a 5-fold cross-validation strategy to validate results over different portions of the input data.

### 4.2 Performance on Voice Command Dataset

Subsequently, performance of the model was tested in a test set that involved voice commands in diverse situations of call for lighting control, temperature regulations, security activation or operation of the entertainment system. The set of data enhanced with noise and various accents allowed for the creation of real-life environments for the smart home's assistants.

### 4.4 Generalization and Robustness

The performance of the model to generalize to unseen commands was fully tested by measuring its performance on set of unseen commands obtained from a test set. The episodic training setup made it possible for even with only a few pieces of new classes, the model could map in new voice command to their corresponding prototypes nicely. To explore its robustness further, we exposed the model to diverse noise conditions, including background chatter, sounds made by household appliances and outdoor noise respectively. The model retained high recognition rate indicating its noise tolerance. The pitch shifting methods and the addition of background noise allowed the model to be flexible to different nuances in the audio input and hence more robust in actual environments.

### 4.5 Model Efficiency and Latency

Real time performance was tested based on model inference run on an edge device, a Raspberry Pi 4. The system was successful to process and classify voice commands with average latency of 30 milliseconds and thus could be deployed at smart homes where fast responses are important.

**Table 3.** Shows the inference time for different command categories

Command Category	Inference Time (ms)
Lighting Control	28
Temperature Adjustment	32
Security Activation	29
Entertainment Control	30

The actual-time classifying feature, combined with the low cost of computation, guarantees that the presented system can be implemented in devices with limited resources for smart homes.

#### 4.6 Limitations and Future Work

Although in controlled environments, the model was successful in recognizing voice commands, difficulties for noisy, real-world scenarios remain. Variations in future works may involve robustness in very dynamic environments where the background noise becomes immeasurable unpredictable and complex. Moreover, the performance of the model in very small datasets could be improved further by investigating other forms of advanced transfer learning.

In addition, the system's scalability to support multi-user situations as well as multi-lingual support would increase its applicability range. Adding speaker identification, or context awareness, could improve the performance of the system in personalized smart homes.

#### 5. CONCLUSION

In this work, we proposed a few-shot learning-based model for voice command recognition in smart homes with the aim of increasing the performance of smart assistants in situations where data is scarce. Using episodic training and few-shot learning, the model was able to accurately recognize voice commands with little labeled data making it suitable in the real world smart home setting where labeled data are often scarce. The model showed high accuracy (94.5%) and good performance in various types of commands: lighting, temperature, security, and entertainment, with an F1-score of 92.9%. It also had wonderful generalization to new unseen commands and did well in noisier levels as well and with the help of data augmentation methods such as pitch shifting and adding noise. In addition, its real time inference capability with mean latency of 30 milliseconds on a Raspberry Pi 4 also confirms its appropriateness for edge deployment in smart home assistants. However, areas of improvement involve increasing robustness in dynamic noise environments, dealing with multi-user situations and dealing with multilingual capabilities. In future work, we are interested in doing transfer learning for multilingual and cross – accent support, and more advanced noise filtering techniques. Overall, the few shot learning approach is found to be a

viable solution to voice command recognition, given a limited amount of resources, and has considerable applicability in the case of real time applications in the IoT system, smart homes, and edge computing devices.

#### REFERENCES

1. Chen, Y., Yang, W., & Xu, C. (2019). Meta-learning for speaker verification: A few-shot approach. *IEEE Transactions on Audio, Speech, and Language Processing*, 27(4), 907-919. <https://doi.org/10.1109/TASLP.2019.2908356>
2. Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., & Weyand, T. (2017). MobileNets: Efficient convolutional neural networks for mobile vision applications. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, 1-9. <https://doi.org/10.1109/CVPR.2017.262>
3. Koch, G., Zemel, R., & Salakhutdinov, R. (2015). Siamese neural networks for one-shot image recognition. *International Conference on Machine Learning (ICML)*, 1-9. <https://arxiv.org/abs/1503.03832>
4. Kim, J., Yoon, H., & Lee, J. (2019). Prototypical networks for spoken keyword recognition. *IEEE Transactions on Neural Networks and Learning Systems*, 30(5), 1381-1392. <https://doi.org/10.1109/TNNLS.2018.2872614>
5. Snell, J., Swersky, K., & Zemel, R. (2017). Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 4080-4090. <https://arxiv.org/abs/1703.05175>
6. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., & Kavukcuoglu, K. (2016). Matching networks for one-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 3630-3638. <https://arxiv.org/abs/1606.04080>
7. Warden, P. (2018). Speech Commands: A dataset for limited-vocabulary speech recognition. *TensorFlow Blog*. [https://www.tensorflow.org/datasets/community\\_catalog/huggingface/speech\\_commands](https://www.tensorflow.org/datasets/community_catalog/huggingface/speech_commands)
8. Zhang, Y., & Lin, Y. (2021). Tiny speaker identification using few-shot learning for microcontrollers. *IEEE Access*, 9, 32534-

32546.  
<https://doi.org/10.1109/ACCESS.2021.3054724>
9. Zhang, Z., Huo, J., Li, Y., & Xue, W. (2017). Hello Edge: Real-time keyword spotting on embedded systems. *Proceedings of the 26th International Conference on Computer Communication and Networks (ICCCN)*, 1-9. <https://doi.org/10.1109/ICCCN.2017.8038437>