



The 39th Annual AAAI  
Conference on Artificial  
Intelligence

FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA



西安交通大学  
XI'AN JIAOTONG UNIVERSITY



兰州大学  
LANZHOU UNIVERSITY



SMILES LAB  
域 览 九 州

极 韵 | DREAM

# M3Net: Efficient Time-Frequency Integration Network with Mirror Attention for Audio Classification on Edge

**Presenter: Xuanming Jiang**

**2025 | Pennsylvania, USA**

**Xuanming Jiang<sup>1,3</sup>, Baoyi An<sup>2,3</sup>, Guoshuai Zhao<sup>1,4\*</sup>, Xueming Qian<sup>1,4</sup>**

<sup>1</sup>School of Software Engineering, Xi'an Jiaotong University, Xi'an, China <sup>2</sup>School of Physical Science and Technology, Lanzhou University, Lanzhou, China

<sup>3</sup>Xi'an Jiyun Technology Co., Ltd., Xi'an, China <sup>4</sup>Shaanxi Yulan Jiuzhou Intelligent Optoelectronic Technology Co., Ltd., Xi'an, China

jiangxm24@stu.xjtu.edu.cn, anby20@lzu.edu.cn, {guoshuai.zhao, qianxm}@mail.xjtu.edu.cn



# Outline

1

**Introduction**

2

**Methodology**

3

**Experimental Results**

4

**Conclusion & Discussion**



The 39th Annual AAAI  
Conference on Artificial  
Intelligence

FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA

1

# Introduction



## Applications:

- Human-machine interaction
- Intelligent robotic
- Wise information technology of med

## Requirement

### High Real-Time Performance

Move the data processing procedure from the cloud to the edge for reducing latency

## Challenges :

- Place higher demands on model performance
- Impose greater computation cost on MCUs with limited resources

## Solution

Apply CNNs for addressing multiple audio classification scenarios based on edge devices

**Higher Model Complexity**

**Generally**

**Better Performance**

## Challenges :

- SOTA methods currently struggle to achieve a required **balance between performance and complexity**
- CNNs are originally designed for image-based tasks, which may limit the performance of audio classification models based on CNNs

## Existing Methods:

- Model compression and pruning
- Multiple attention mechanisms
- Transfer learning



The 39th Annual AAAI  
Conference on Artificial  
Intelligence

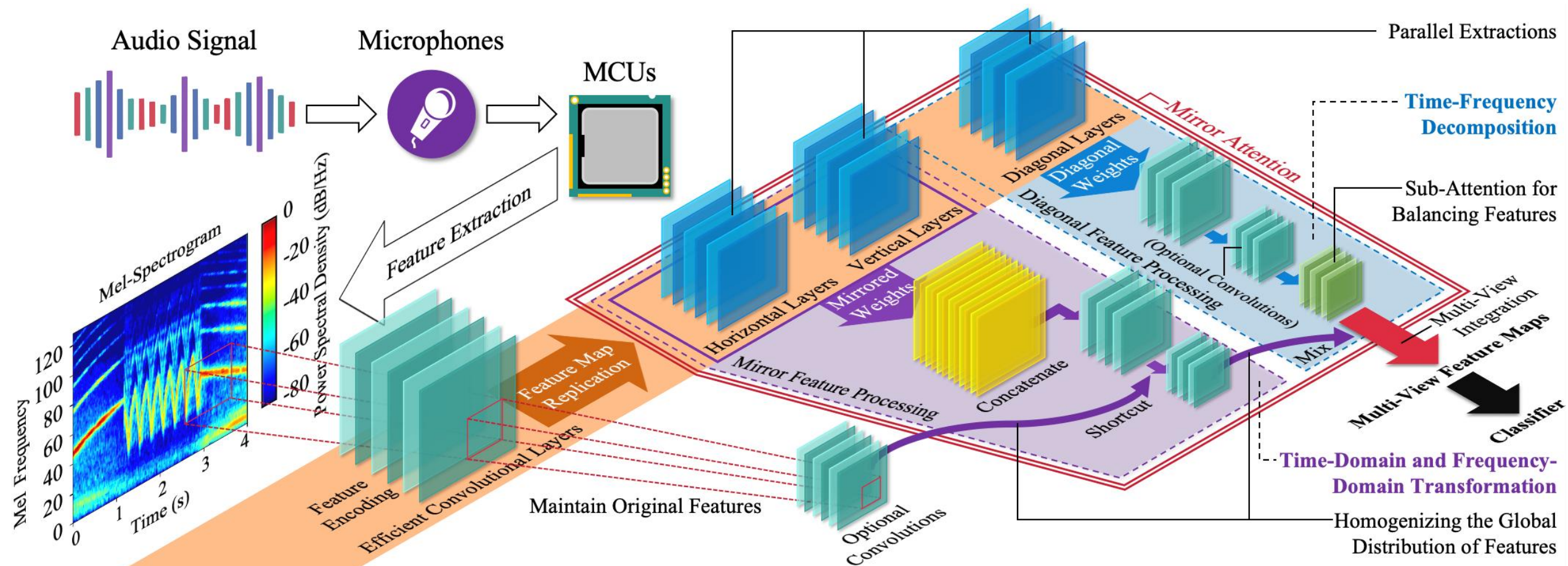
FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA

2

# Methodology



## Mini Mirror Multi-View Network (M3Net)



**(1) Define feature element:**

$$\mathbf{X} = \{X(w, h, c) | w \in [0, W-1], h \in [0, H-1], c \in [0, C-1]\}$$

**Mirror Attention**

$$\mathbf{X}_{MA}(w, h, c) = \mathbf{X}_{MFP}(w, h, c) \odot \mathbf{X}_{DFP}(w, h, c)$$

**(2) Time-domain transformation:**

$$X_{T,c \in [0,C-1]}(w, h) = M_{HMF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{HMF} = \begin{bmatrix} -1 & 0 & W \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

**Mirror features between time and frequency****(3) Frequency-domain transformation:**

$$X_{F,c \in [0,C-1]}(w, h) = M_{VMF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{VMF} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & H \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{X}_{MFP} = F_M(M_{HMF}\mathbf{X}, M_{VMF}\mathbf{X}), \quad \mathbf{X} \in \mathbb{R}^{W \times H \times C}$$

**(4) Time-frequency decomposition:**

$$X_{D,c \in [0,C-1]}(w, h) = M_{CSF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{CSF} = \begin{bmatrix} -1 & 0 & W \\ 0 & -1 & H \\ 0 & 0 & 1 \end{bmatrix}$$

**Separated time and frequency features**

$$\rightarrow \mathbf{X}_{DFP} = F_D(M_{CSF}\mathbf{X}) = F_D(M_{HMF}M_{VMF}\mathbf{X})$$



## (1) Define feature element:

$$\mathbf{X} = \{X(w, h, c) | w \in [0, W-1], h \in [0, H-1], c \in [0, C-1]\}$$

**Mirror Attention**

$$\mathbf{X}_{MA}(w, h, c) = \mathbf{X}_{MFP}(w, h, c) \odot \mathbf{X}_{DFP}(w, h, c)$$

## (2) Time-domain transformation:

$$X_{T,c \in [0,C-1]}(w, h) = M_{HMF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{HMF} = \begin{bmatrix} 1 & 0 & W \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

**Mirror features between time and frequency**

## (3) Frequency-domain transformation:

$$X_{F,c \in [0,C-1]}(w, h) = M_{VMF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{VMF} = \begin{bmatrix} 0 & 0 & W \\ 0 & -1 & H \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{X}_{MFP} = F_M(M_{HMF}\mathbf{X}, M_{VMF}\mathbf{X}), \quad \mathbf{X} \in \mathbb{R}^{W \times H \times C}$$

## (4) Time-frequency decomposition:

$$X_{D,c \in [0,C-1]}(w, h) = M_{CSF} \begin{bmatrix} w \\ h \\ \theta \end{bmatrix}, \quad M_{CSF} = \begin{bmatrix} 0 & 0 & W \\ 0 & -1 & H \\ 0 & 0 & 1 \end{bmatrix}$$

**Separated time and frequency features**

$$\mathbf{X}_{DFP} = F_D(M_{CSF}\mathbf{X}) = F_D(M_{HMF}M_{VMF}\mathbf{X})$$





The 39th Annual AAAI  
Conference on Artificial  
Intelligence

FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA

3

# Experimental Results

## UrbanSound 8K

**Volume:** ~8,700 audio clips**Content:** Typical urban noises**Duration:** <4s

**Details:** 10 classes in total, including air conditioners, car horns, children playing, dogs barking, drilling, engine idling, gunshots, jackhammers, sirens, and street music.

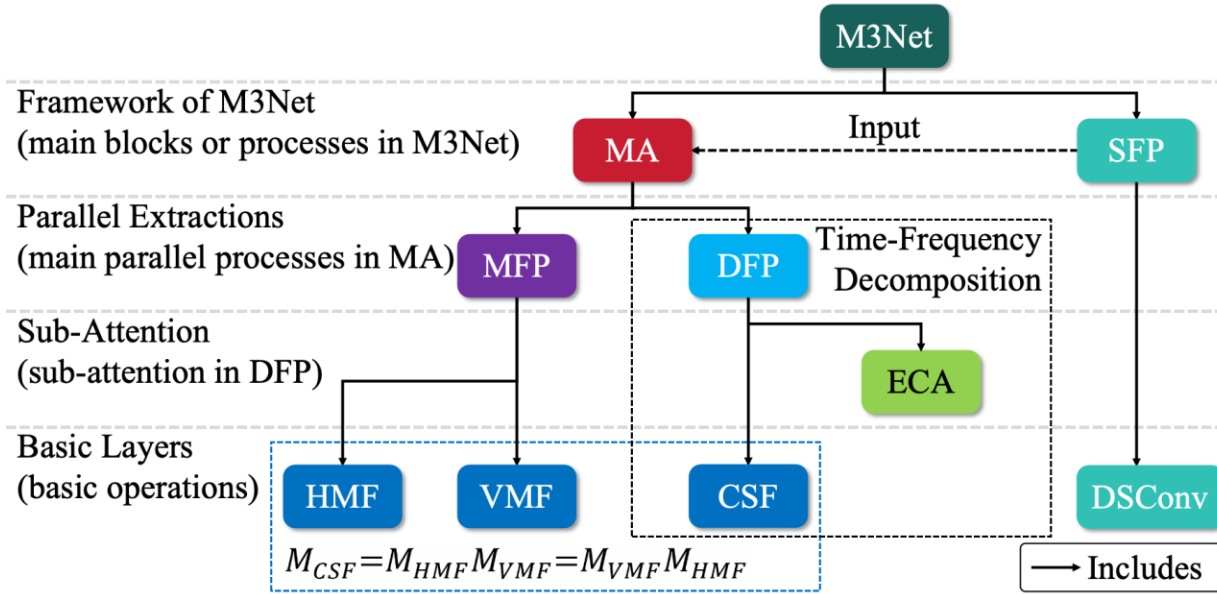
Dataset: UrbanSound8K		Accuracy	# Param	Wilcoxon	Data	Pre-training
Methods	Features	(%)	( $\times 10^6$ )	p-value	augmentation	
AemNet-DW	Log-Mel	82.25	0.9	$< 5.0e-2$		✓
ULSED	Log-Mel	83.5	0.34	$< 5.0e-2$	✓	
2D CNN	GFCC	89	1.8	$< 5.0e-2$	✓	
PhiNets M40	Mel-Spectrogram	76.3	<b>0.027</b>	$< 5.0e-2$	✓	
SE-TCAM 1D CNN	Raw audio	94.04	0.81	$< 5.0e-2$		
<b>M3Net (Ours)</b>	Mel-Spectrogram	<b>97.44</b>	0.029	baseline		

## SpeechCommands V2

**Volume:** ~105,000 audio clips**Content:** Voice commands**Duration:** ~1s

**Details:** 35 classes in total, including verbal directives such as “yes”, “no”, “up”, “down”, numeric commands ranging from “zero” to “nine”, miscellaneous terms like “bed”, “bird”, “cat”, “dog”.

Dataset: SpeechCommandsV2		Accuracy	# Param	Wilcoxon	Data	Pre-training
Methods	Features	(%)	( $\times 10^6$ )	p-value	augmentation	
DeLoRes M	Log-Mel	89.7	5.3	$< 5.0e-2$	✓	✓
AdaptFormer	Log-Mel	92.3	1.43	$< 5.0e-2$		✓
DCLS-Delays	Mel-Spectrogram	95.35	2.5	$< 5.0e-2$	✓	
SeqBoat	Raw audio	97.35	0.293	$2.2e-1$		
EAT-S-GMME	Raw audio	<b>97.88</b>	1.54	$1.0e-1$	✓	
<b>M3Net (Ours)</b>	Mel-Spectrogram	97.03	<b>0.031</b>	baseline		



Dataset	M3Net w/o MA		M3Net w/o SFP	
	Accuracy (%)	Wilcoxon	Accuracy (%)	Wilcoxon
US8K	75.29 (-22.15)	< 5.0e-2	87.31 (-10.13)	< 5.0e-2
SCV2	68.13 (-28.90)	< 5.0e-2	84.26 (-12.77)	< 5.0e-2

Ablation target	1	2	3	4
ECA		✓	✓	✓
CSF			✓	✓
VMF				✓
HMF				
US8K	Accuracy	93.88	94.56	95.36
	(%)	(-3.56)	(-2.88)	(-2.08)
SCV2	Accuracy	94.00	94.20	94.89
	(%)	(-3.03)	(-2.76)	(-2.14)
US8K	Wilcoxon	< 5.0e-2	< 5.0e-2	< 5.0e-2
SCV2	Wilcoxon	< 5.0e-2	< 5.0e-2	< 5.0e-2



The 39th Annual AAAI  
Conference on Artificial  
Intelligence

FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA

4

# Conclusion & Discussion

- ① With the advancement of AI, edge devices are taking on increasing responsibilities
- ② As IC enter the post-Moore's Law era, edge devices' performance limitations become increasingly significant
- ③ Edge devices need to meet societal demands from both hardware and algorithmic perspectives

### Background

- ① An **time-frequency separation** method for audio learning
- ② An **mirror attention mechanism** for extending pattern set
- ③ The **M3Net** with high-performance but lightweight

### Contributions

M3Net transcends the conventional emulation of human auditory perception by CNNs and ventures into the field of replicating human counterintuitive supervised cognitive processes.

### Methodology

M3Net:

- ① Achieving comparable performance to SOTAs with **less than 10%** of the parameters
- ② Adapting to different kinds of audio contents
- ③ More sensitive to complex and dynamic audio

### Achievements



The 39th Annual AAAI  
Conference on Artificial  
Intelligence

FEBRUARY 25 – MARCH 4, 2025 |  
PHILADELPHIA, PENNSYLVANIA, USA



西安交通大学  
XI'AN JIAOTONG UNIVERSITY



极韵 | DREAM



兰州大学  
LANZHOU UNIVERSITY



SMILES LAB  
域 览 九 州

# Thank you

2025 | PHILADELPHIA, USA