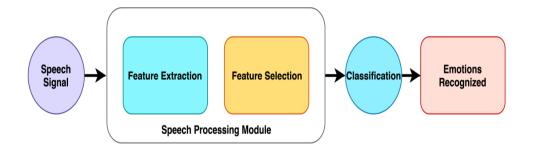
BAB II

TINJAUAN PUSTAKA

2.1 Landasan Teori

2.1.1 Speech Emotion Recognition (SER)

Speech Emotion Recognition (SER) dapat didefinisikan sebagai suatu sistem pengenalan pola. Hal ini menunjukkan bahwa tahapan-tahapan yang terdapat pada sistem pengenalan pola juga terdapat pada sistem pengenalan emosi berbasis ucapan (Selvaraj dkk., 2016). Gambar 2.1 merupakan alur dari sistem yang disederhanakan yang digunakan untuk Speech Emotion Recognition.



Gambar 2.1 Tahapan SER

Berdasarkan Gambar 2.1, *Speech Emotion Recognition* berisi lima modul utama yaitu, input ucapan emosional, ekstraksi fitur, pemilihan fitur, klasifikasi, dan output emosional yang dikenali (El Ayadi dkk., 2011).

Proses untuk mengenali sekumpulan emosi secara otomatis menjadi perhatian utama dalam SER. Oleh karena itu untuk mengklasifikasikan sejumlah besar data emosi sangatlah rumit. Masalah SER telah dibahas selama beberapa tahun menggunakan metode statistik dan algoritma *machine*

learning, seperti *Support Vector Machines* (SVMs) dan berbagai algoritma regresi (Tarantino dkk., 2019).

2.1.2 Klasifikasi

Klasifikasi adalah proses menemukan suatu model yang menggambarkan dan membedakan kelas data atau konsep, dengan tujuan agar model tersebut dapat digunakan untuk memprediksi kelas dari objek yang label kelasnya belum diketahui. Selain itu, klasifikasi digunakan secara luas pada *data mining* untuk mengelompokkan data ke dalam beragam kelas (Hasanain & Rizal, 2021). Teknik dari klasifikasi adalah dengan melihat variabel dari kelompok data yang sudah ada. Klasifikasi bertujuan untuk memprediksi kelas dari suatu objek yang tidak diketahui sebelumnya. Umumnya klasifikasi terdiri dari tiga tahap, yaitu pembangunan model, penerapan model, dan evaluasi (Nasution dkk., 2019).

2.1.3 Convolutional Neural Network (CNN)

CNN adalah pengembangan dari *Multilayer Perceptron* (MLP) yang didesain untuk mengolah data dua dimensi. CNN termasuk dalam jenis *Deep Neural Network* karena kedalaman jaringan yang tinggi dan banyak diaplikasikan pada data data dua dimensi (Nurona Cahya dkk., 2021). Cara kerja CNN memiliki kesamaan pada MLP, namun dalam CNN setiap neuron dipresentasikan dalam bentuk dua dimensi, tidak seperti MLP yang setiap neuron hanya berukuran satu dimensi (Hasanain & Rizal, 2021).

Menurut (Martiyaningsih, 2022) secara umum CNN memiliki beberapa layer, yaitu *Convolutional Layer, Pooling Layer*, dan *Fully Connected Layer*.

- Convolutional Layer, merupakan tahap dimana seluruh data menyentuh lapisan convolutional yang mengalami proses konvolusi dan akan difilter kemudian menghasilkan sebuah activation map. Lapisan pada layer ini memiliki tiga parameter, yaitu depth, stride, dan zero padding.
- 2. *Pooling Layer*, merupakan lapisan yang menggunakan *Feature Map* sebagai input dan mengolahnya berdasar nilai piksel terdekat. Lapisan ini disisipkan ke dalam lapisan konvolusi secara teratur. Bentuk Pooling Layer yang paling umum adalah 2x2.
- Fully Connected Layer, merupakan lapisan yang digunakan untuk tujuan melakukan transformasi pada dimensi data agar dapat diklasifikasikan secara linear.

2.1.4 Mel Frequency Cepstral Coefficient (MFCC)

Mel Frequency Cepstral Coefficient (MFCC) merupakan metode ekstraksi fitur untuk mendapatkan lebih banyak informasi dalam data audio. MFCC dapat digunakan untuk mengekstrak data sinyal suara sehingga bisa didapatkan ciri yang terdapat pada data sinyal suara (Emanuella dkk., 2021). MFCC bekerja menyerupai adaptasi sistem pendengaran manusia, dengan cara memfilter sinyal suara secara logaritmik untuk frekuensi tinggi yang nilainya di atas 1000 Hz dan memfilter sinyal suara secara linear untuk frekuensi rendah yang nilainya di bawah 1000 Hz (Danika dkk., 2023). Metode ini efektif dalam menangkap karakteristik penting dari suara karena didasarkan pada persepsi pendengaran manusia.

Proses ekstraksi fitur MFCC terdiri dari beberapa tahapan penting. Dimulai dengan *pre-emphasis* untuk meningkatkan energi sinyal pada frekuensi tinggi, dilanjutkan dengan *framing* yang membagi sinyal menjadi frame-frame pendek. Selanjutnya, *windowing* dilakukan untuk mengurangi diskontinuitas di tepi *frame*. *Fast Fourier Transform* (FFT) kemudian digunakan untuk mengubah sinyal dari domain waktu ke domain frekuensi (Gupta and Gupta, 2016; Ralev and Krastev, 2023).

Setelah transformasi ke domain frekuensi, spektrum diproses melalui serangkaian filter segitiga yang terdistribusi secara mel, yang dikenal sebagai *Mel Filter Bank*. Langkah berikutnya adalah *logarithmic compression* untuk menyesuaikan dengan persepsi manusia terhadap loudness. Kemudian, *Discrete Cosine Transform* (DCT) digunakan untuk mengubah log mel spectrum menjadi *domain cepstral* (Sabitha V, 2013; Gupta and Gupta, 2016).

MFCC telah terbukti efektif dalam berbagai aplikasi pemrosesan suara. Dalam konteks pengenalan emosi, MFCC mampu menangkap perubahan dalam nada suara dan intonasi yang sering dikaitkan dengan ekspresi emosional. Penelitian menunjukkan bahwa MFCC, ketika dikombinasikan dengan algoritma pembelajaran mesin seperti *Support Vector Machine* (SVM) atau *Convolutional Neural Network* (CNN), dapat mencapai akurasi yang tinggi dalam klasifikasi emosi (Ravi and Taran, 2024). Meskipun MFCC sangat efektif, beberapa penelitian telah mengusulkan modifikasi atau alternatif untuk meningkatkan kinerjanya. Misalnya, *Power-Normalized Cepstral*

Coefficients (PNCC) dan Modified Group Delay Function (ModGDF) telah diusulkan sebagai alternatif yang potensial, terutama dalam kondisi berderau.

Dalam implementasinya, MFCC telah diterapkan dalam berbagai domain di luar pemrosesan suara manusia. Termasuk klasifikasi suara lingkungan, deteksi kebocoran dalam jaringan distribusi air, dan bahkan dalam analisis getaran untuk diagnosis kesalahan mekanis. Keunggulan MFCC terletak pada kemampuannya untuk menangkap karakteristik spektral suara yang relevan dengan persepsi manusia (Razak et al., 2008; das, Jena and Barik, 2014).

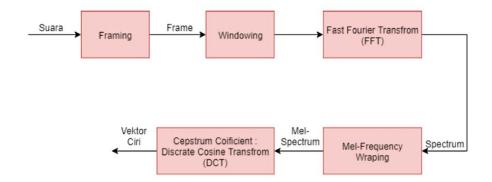
Dengan perkembangan teknologi pembelajaran mendalam, MFCC sering digunakan sebagai input untuk arsitektur *neural network* yang kompleks. Misalnya, kombinasi MFCC dengan *Long Short-Term Memory* (LSTM) *networks* telah menunjukkan hasil yang menjanjikan dalam tugas-tugas seperti pengenalan emosi. Selain itu, penggunaan MFCC bersama dengan teknik augmentasi data seperti *gaussian noise, time stretching*, dan *pitch shifting* telah terbukti dapat meningkatkan kinerja model (Padman and Magare, 2023; Neili and Sundaraj, 2024).

Dalam konteks pengenalan emosi dari suara, MFCC telah menunjukkan keefektifannya ketika digabungkan dengan berbagai teknik pembelajaran mesin. Penelitian terbaru menunjukkan bahwa penggunaan MFCC bersama dengan model seperti BiLSTM atau CNN dapat mencapai akurasi yang tinggi dalam klasifikasi emosi. Beberapa studi bahkan menggabungkan MFCC dengan fitur tekstual dan visual untuk meningkatkan akurasi klasifikasi emosi

dalam konteks multimodal (Mishra, Warule and Deb, 2024; Nikhila Gudipati et al., 2024).

MFCC juga telah diaplikasikan dalam berbagai bidang lain seperti klasifikasi *genre* musik, deteksi suara sintetis, dan bahkan dalam pemantauan kesehatan. Dalam klasifikasi genre musik, MFCC telah digunakan untuk mengekstrak fitur dari lagu-lagu dalam berbagai bahasa, termasuk Bengali. Sementara itu, dalam konteks kesehatan, MFCC telah digunakan untuk menganalisis suara batuk dan pola pernapasan untuk diagnosis penyakit pernapasan (Khan and Hafiz, 2024; Sai Varun et al., 2024).

Meskipun MFCC telah terbukti sangat efektif, penelitian terus berlanjut untuk meningkatkan kinerjanya dan mengeksplorasi kombinasinya dengan teknik-teknik baru. Beberapa penelitian terbaru fokus pada optimalisasi parameter MFCC, penggabungannya dengan fitur lain seperti spektrogram, dan penggunaan teknik pembelajaran mesin yang lebih canggih untuk meningkatkan akurasi klasifikasi. Dengan perkembangan ini, MFCC terus menjadi salah satu metode ekstraksi fitur yang paling penting dan banyak digunakan dalam pemrosesan sinyal suara (Zhong, 2023; Bai, 2024; Gourisaria et al., 2024).



Gambar 2.2 Alur Ekstraksi Fitur MFCC

Gambar 2.2 menjelaskan alur dari tahapan ekstraksi fitur MFCC. Pertama, sinyal suara dipotong untuk menghilangkan keheningan atau gangguan yang mungkin muncul pada awal maupun akhir suara dilakukan framing, membagi ke dalam sejumlah frame. Selanjutnya, proses windowing digunakan untuk meminimalkan diskontinuitas sinyal. Fast Fourier Transform (FFT) kemudian diterapkan untuk mengubah setiap frame ke domain frekuensi. Dalam Mel-frequency wrapping block, sinyal diplot terhadap spektrum Mel untuk meniru pendengaran manusia. Terakhir, Discrete Cosine Transform (DCT) dilakukan untuk menghasilkan vektor ciri (Handoko & Suyanto, 2019).

2.1.5 Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D)

Crowd-sourced Emotional Multimodal Actors Dataset (CREMA-D) adalah dataset berlabel yang digunakan untuk mempelajari pengenalan ekspresi dan emosi. CREMA-D mencakup 7.442 klip dari 91 aktor dan aktris dengan usia dan etnis yang beragam dengan mengekspresikan enam bentuk emosi secara universal yaitu anger, disgust, fear, happy, neutral and sad (Cao dkk., 2014).

Secara umum, usia dan latar belakang etnis dapat memengaruhi cara seseorang berbicara, termasuk pola naik-turun nada suaranya (intonasi). Menurut Jacewicz dan Fox (2013), perbedaan usia memengaruhi anatomi dan fungsi pita suara, misalnya kelenturan dan ketebalan lipatan vokal, sehingga nada dasar serta rentang frekuensi suara bisa berbeda pada tiap kelompok umur. Di sisi lain, latar belakang etnis berkaitan dengan dialek atau aksen tertentu, yang juga membentuk ciri khas intonasi setiap kelompok budaya (Laver, 1994). Dengan demikian, kombinasi antara faktor usia dan etnis dapat menghasilkan karakter suara yang beragam, termasuk pola intonasi yang berbeda.

Dalam SER, latar belakang etnis dan dialek berhubungan erat dengan variasi cara berbicara, seperti perbedaan aksen, pilihan kata, dan pola naikturun nada (intonasi). Faktor-faktor ini dapat memengaruhi ciri akustik yang diandalkan oleh sistem SER, misalnya frekuensi dasar suara (pitch) atau kecepatan berbicara, sehingga model yang tidak dilatih dengan data dari berbagai etnis dan dialek berisiko kesulitan mengenali emosi pada penutur tertentu. Dalam salah satu penelitian ditunjukkan bahwa perbedaan aksen dapat membuat pola suara marah atau senang terdengar berbeda, sehingga model SER yang hanya berfokus pada satu dialek atau aksen tertentu cenderung menghasilkan akurasi yang lebih rendah bila diterapkan pada penutur dengan latar belakang bahasa atau etnis lain (Elbarougy et al., 2014). Dengan demikian, variasi etnis dan dialek penting diperhatikan untuk memastikan bahwa sistem SER mampu mengakomodasi seluruh pengguna secara andal.

2.2 Penelitian Terkait dan Kebaruan Penelitian

2.2.1 State of The Art

Tabel 2.1 menunjukan perbandingan penelitian yang berhubungan dengan fokus pada permasalan penelitian dan hasil penelitian yang dihasilkan dalam proses klasifikasi.

Tabel 2.1 State of The Art

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
1	Rizqi Fathin Fadhillah,	2023	Klasifikasi Suara Untuk	Beberapa upaya untuk	Berdasarkan hasil penelitian
	Raden Sumiharto		Memonitori Hutan Berbasis	mencegah aktivitas ilegal di	ini, didapatkan bahwa model
			Convolutional Neural	hutan telah dilakukan seperti	MobileNetV3-Small + MLP
			Network	patroli	dengan data latih gabungan
				oleh petugas setempat. Akan	dari augmentasi time stretch
				tetapi upaya tesebut	dan time shift memberikan
				mengalami beberapa	performa yang paling bagus

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
				Irandala diantamanya adalah	untuk diimalamantasikan
				kendala, diantaranya adalah	untuk diimplementasikan
				sumber daya manusia yang	pada sistem ini, dengan
				kurang, sarana dan prasarana	durasi inference 0.8 detk;
				yang kurang mendukung,	akurasi sebesar 93.96%; dan
				dan keterbatasan dana.	presisi sebesar 94.1%.
				Patroli yang dilakukan oleh	
				manusia juga terbatas pada	
				waktu tertentu, sehingga	
				pada waktu yang lain masih	
				berpotensi adanya tindakan	
				ilegal.	

No	Nama Peneliti	Tahun	Judu	I	Masalah		F	Iasil
2	Andani Achmad,	2022	Klasifikasi	Penyakit	Saat ini, sistem ke	esehatan	Hasil	pengujian
	Adnan, Muhammad		Pernapasan	Berbasis	pernafasan	telah	mengemukal	kan bahwa
	Rijal		Visualisasi	Suara	memanfaatkan	metode	metode Su	ipport Vector
			Menggunakan	Metode	machine learning	untuk	Machine ber	hasil diterapkan
			Support Vector N	Machine	melakukan kl	lasifikasi	dalam	mengklasifikasi
					penyakit pernapasan	n melalui	penyakit asm	na, bronkitis dan
					visualisasi suara s	stetoskop	tuberkulosis	dengan akurasi
					yang merupakan ir	nformasi	sebesar 46.3°	7%.
					grafik bunyi dalam	rentang		
					waktu dan f	frekuensi		
					tertentu.			

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
3	Cecilia Tania	2021	Klasifikasi Suara Kucing	Kucing dan anjing	Metode praproses dan
	Emanuella, Musfita, dan		dan Anjing Menggunakan		metode klasifikasi dengan
	Armin Lawi		Convolutional Neural	peliharaan manusia yang	menggunakan
			Network	berkomunikasi melalui	Convolutional Neural
				suara. Untuk membantu	Network cukup baik untuk
				mengetahui, membedakan	menentukan kebenaran dari
				dan mengenali suara hewan	klasifikasi data audio. Hal
				dengan peliharaan yang lain,	ini terbukti dengan hasil
				dibutuhkan klasifikasi suara	akurasi sebesar 88%.
				hewan agar lebih jelas.	Perubahan tingkat confusion
					tidak mempengaruhi hasil
					akurasi. Hal ini

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
					membuktikan bahwa klasifikasi menggunakan
					metode CNN relatif baik
					terhadap perubahan parameter yang dilakukan.
4	ABHISHEK SEHGAL	2021	A Convolutional Neural	VAD digunakan untuk	Makalah ini telah
	AND NASSER		Network Smartphone App	tujuan di mana klasifikasi	menyediakan jaringan saraf
	KEHTARNAVAZ		for Real-Time Voice	atau estimasi kebisingan	tiruan convolutional untuk
			Activity Detection	diaktifkan oleh VAD untuk	melakukan deteksi aktivitas
				menyesuaikan parameter	suara secara realtime dengan
				algoritma pengurangan	latensi audio yang rendah.
				derau tergantung pada kelas	Aplikasi ini telah

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
				atau jenis derau. Untuk	dikembangkan untuk
				bagian sinyal atau bagian di	smartphone Android dan
				mana ucapan dalam derau	iOS. Arsitektur dari jaringan
				atau ucapan+derau	saraf convolutional telah
				terdeteksi, tidak ada	dioptimalkan untuk
				klasifikasi/estimasi derau	memungkinkan bingkai
				yang dilakukan dan	audio untuk diproses secara
				pengurangan derau	real-time tanpa bingkai apa
				dilakukan berdasarkan jenis	pun yang dilewati dengan
				derau yang diidentifikasi	tetap mempertahankan
				terakhir.	akurasi suara yang tinggi
					deteksi aktivitas.

No	Nama Peneliti	Tahun	J	udul			Masalah			Hasil
5	Muhammad Hasbi	2020	Klasifikasi	Suara	Paru	Deteksi k	elainan p	ada suara	Setelah me	lakukan pelatihan
	Ashshiddieqy, Jondri,		Dengan	Convolu	tional	paru biasa	a dilakuka	an dengan	pada mod	lel pembelajaran
	Achmad Rizal		Neural Netw	ork (CNN)	mengguna	akan s	stetoskop.	mesin	CNN dengan
						Pendeteks	sian sua	ara paru	menggunal	kan generalisasi
						dengan s	tetoskop	memiliki	berupa a	nugmentasi dan
						keterbatas	san karer	na sangat	dropout	layer, maka
						bergantun	g pada	individu	didapatkan	akurasi sebesar
						yang	n	nelakukan	84,80% te	rhadap data latih
						pemeriksa	aan.		dan 78,09	% terhadap data
									uji.	
6	Mazin Abed	2020	Voice Patho	ology Dete	ection	Patologi	suara	memiliki	Hasil	eksperimen
	Mohammed, Karrar		and Classi	fication (Using	dampak	negatit	f pada	menunjukk	can metode CNN

No	Nama Peneliti	Tahun	Judul		Masalah	Hasil
	Hameed Abdulkareem,		Convolutional Neur	ral	keteraturan getaran dan	yang diusulkan untuk
	Salama A. Mostafa,		Network Model		fungsi suara, yang mengarah	deteksi patologi ucapan
	Mohd Khanapi Abd				pada peningkatan	mencapai akurasi hingga
	Ghani, Mashael S.				kebisingan suara. Suara	95,41%. Metode ini juga
	Maashi, Begonya				normal berubah menjadi	memperoleh 94,22% dan
	Garcia-Zapirain, Ibon				tegang, lemah dan serak	96,13% untuk F1-Score dan
	Oleagordia, Hosam				yang mempengaruhi	Recall. Sistem yang
	Alhakami and Fahad				kualitas. Sampai saat ini,	diusulkan menunjukkan
	Taha AL-Dhief				metode deteksi patologi	kemampuan yang tinggi dari
					suara yang ada saat ini	aplikasi klinis nyata yang
					memiliki evaluasi yang bias	menawarkan solusi
						diagnosis dan perawatan

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
				berdasarkan hal-hal yang	otomatis yang cepat dalam
				bersifat subjektif.	waktu 3 detik untuk
					mencapai akurasi
					klasifikasi.
7	Andre Danika	2020	SENTIMENT ANALYSIS	Perbedaan utama kedua jenis	Hasil dari proses tersebut
	Jangkung Raharjo		ONLINE SHOP ON THE	senar adalah suara yang	kemudian akan
	Bambang Hidayat		PLAY STORE USING	dihasilkan. Senar baja	diklasifikasikan dengan
			METHOD SUPPORT	cenderung menghasilkan	metode support vector
			VECTOR MACHINE	suara yang lebih nyaring,	machine (SVM) dengan
			(SVM)	sedangkan senar nylon	fungsi kernel RBF sebagai
				cenderung menghasilkan	fungsi terbaik dengan
				suara yang mellow. Selain	akurasi 95%. Gitar senar

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
				itu senar baja juga	baja cenderung
				menghasilkan volume suara	menghasilkan frekuensi
				yang lebih besar	maksimum yang lebih besar
				dibandingkan dengan suara	dibandingkan dengan senar
				yang dihasilkan senar nylon.	nylon.
8	Raditya Budi Handoko	2019	Klasifikasi Gender	Banyak metode telah	Sistem klasifikasi gender
	dan Suyanto		Berdasarkan Suara	diusulkan untuk	berdasarkan suara
			Menggunakan Support	membangun sistem	menggunakan SVM yang
			Vector Machine	klasifikasi gender berbasis	dibangun mampu
				suara yang berakurasi tinggi,	memberikan
				namun umumnya masih	akurasi sebesar 100%.
					Pemilihan parameter dan

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
				kurang tahan terhadap derau	kernel SVM sangat
				Kurang tahan terhadap derad	Reffici 5 v ivi sungut
				atau noise.	mempengaruhi akurasi
					sistem. Kernel RBF dan
					Sigmoid mempunyai akurasi
					lebih buruk dibanding
					Linear dan Polynomial (d
					=1).
9	Aditya Singgi Prayogi,	2019	Klasifikasi Suara Tangisan	Bagi orang tua yang	Akurasi terbaik pada proses
	Transparsing Fragogi,	2019	Tradition Saura Tangisan	Bugi Grang taa yang	rikurusi terourk puda proses
	Maulana Rizqi, Tresna		Bayi Berdasarkan Prosodic	mengasuh bayi tentu sulit	klasifikasi menggunakan
	Maulana Fahrudin		Features Menggunakan	untuk menentukan apa yang	data sampling Percentage
			Metode Moments of	diinginkan oleh bayi, karena	Rate yaitu 76% dimana nilai
				suara tangisan yang	K yang digunakan adalah 9.

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
			Distribution dan K-Nearest	dikeluarkan hampir sama.	Sedangkan akurasi terbaik
			Neighbours	Padahal ada perbedaan arti	pada proses klasifikasi
				tertentu dari suara tangisan	menggunakan data sampling
				bayi tersebut.	Leave One Out yaitu 42%
					dengan nilai K yang
					digunakan adalah 5.
10	Nur Hudha Wijaya,	2017	KLASIFIKASI SUARA	Untuk melakukan diagnose	Pendekatan ciri statistis
	Indah Soesanti, Eka		JANTUNG	suara jantuk normal atau	dengan menghitunng nilai
	Firmansyah		MENGGUNAKAN	abnormal (disebut murmur	mean, mode, variance,
			NEURAL NETWORK	patologis) diperlukan	deviation, skewness,
			BACKPROPAGATION	kepekaan dan pengalaman	kurtosis, entropy klasifikasi
				oleh dokter, dengan	dengan neural

No	Nama Peneliti	Tahun	Judul	Masalah	Hasil
No	Nama Peneliti	Tahun	Judul BERBASIS CIRI STATISTIS	Masalah demikian hasil diagnose sangat dipengaruhi oleh subjektivitas dokter.	backpropagation
					backpropagation menunjukkan accuracy mencapai 91,72%.

Penelitian-penelitian sebelumnya dalam bidang pengenalan emosi berbasis suara telah menggunakan berbagai metode klasifikasi, seperti *Hidden* Markov Model (HMM), Gaussian Mixture Model (GMM), Support Vector Machine (SVM), dan metode deep learning, terutama Convolutional Neural Network (CNN) (Achmad dkk., 2022; Ashshiddiegy dkk., 2020; Danika dkk., 2023; Emanuella dkk., 2021; Fadhillah & Sumiharto, 2023; Handoko & Suyanto, 2019; Mohammed dkk., 2020; Prayogi dkk., 2019; Sehgal & Kehtarnavaz, 2018; Wijaya dkk., 2017). Meskipun model-model tradisional seperti HMM dan SVM telah digunakan dengan cukup baik dalam beberapa aplikasi, akurasi yang dihasilkan masih terbatas, khususnya ketika menghadapi dataset yang tidak seimbang dan heterogen. Sebaliknya, penggunaan CNN dalam penelitian lebih mutakhir menunjukkan hasil yang lebih menjanjikan, seperti pada penelitian menggunakan dataset SAVEE yang mencapai akurasi 88% pada data latih. Namun, model-model ini masih menghadapi masalah overfitting ketika diterapkan pada data uji, dengan akurasi yang menurun hingga 52%, yang menunjukkan kebutuhan akan pendekatan yang lebih robust.

Selain itu, banyak penelitian yang menggunakan teknik augmentasi data seperti *gaussian noise, time stretching*, dan *pitch shifting* untuk memperkaya variasi data latih dan meningkatkan kinerja model. Teknik-teknik ini telah terbukti efektif dalam mengatasi masalah ketidakseimbangan kelas emosi dalam dataset, serta meningkatkan kemampuan model untuk mengenali emosi dalam berbagai kondisi suara. Meskipun demikian, beberapa penelitian menunjukkan bahwa pengenalan emosi masih terbatas pada jenis-jenis emosi

yang paling dominan dalam dataset, sementara emosi yang lebih kompleks atau langka sering kali kurang terdeteksi. Oleh karena itu, masih ada tantangan besar dalam meningkatkan generalisasi model agar dapat mengenali berbagai emosi dengan lebih akurat pada data uji yang beragam.

Penelitian ini bertujuan untuk mengisi kesenjangan yang ada dengan menggunakan dataset CREMA-D, yang menawarkan variasi emosi dan intensitas suara yang lebih luas. Dataset ini lebih representatif karena mencakup berbagai variasi nada suara, dan situasi yang relevan untuk pengenalan emosi berbasis suara. Selain itu, penelitian ini mengusulkan penggunaan teknik augmentasi data yang lebih beragam untuk memperkaya data pelatihan, serta peningkatan metode ekstraksi fitur untuk mengatasi masalah overfitting dan ketidakseimbangan kelas dalam dataset. Dengan pendekatan ini, diharapkan penelitian ini dapat meningkatkan kinerja model CNN dalam mengenali emosi pada suara dengan lebih akurat, serta memberikan kontribusi signifikan untuk pengembangan teknologi *Speech Emotion Recognition* (SER) yang lebih efektif dan dapat digeneralisasi ke berbagai aplikasi nyata.

2.2.2 Matriks Penelitian

Matriks Penelitian merupakan matriks yang berisi gambaran keseluruhan isi penelitian terkait dan penelitian yang akan dilakukan. Matriks penelitian pada penelitian ini dapat dilihat pada Tabel 2.2.

Tabel 2.2 Matriks Penelitian

No	Penulis / Tahun	Judul	Metode					Tujuan		
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
1	Rizqi Fathin Fadhillah, Raden Sumiharto	Klasifikasi Suara Untuk Memonitori Hutan Berbasis Convolutional Neural		V				√		
	(2023)	Network								
2	Andani Achmad, Adnan, Muhammad Rijal (2022)	KLASIFIKASI PENYAKIT PERNAPASAN BERBASIS VISUALISASI SUARA			V					

No	Penulis / Tahun	Judul	Metode					Tujuan		
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
		MENGGUNAKAN								
		METODE SUPPORT								
		VECTOR MACHINE								
3	Andre Danika	Deteksi Suara Gitar								
	Jangkung	Dengan Bahan Jenis								
	Raharjo	Senar Berbeda								
	Bambang	Melalui Ciri Akustik								
	Hidayat (2022)	Dengan Mel-			$\sqrt{}$		$\sqrt{}$		$\sqrt{}$	
		Frequency Cepstral								
		Coefficients (MFCC)								
		Dan Support Vector								
		Machine (SVM)								

No	Penulis / Tahun	Judul	Metode					Tujuan		
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
4	Cecilia Tania	Klasifikasi Suara								
	Emanuella,	Kucing dan Anjing								
	Musfita, dan	Menggunakan						$\sqrt{}$		
	Armin Lawi	Convolutional Neural								
	(2021)	Network								
5	ABHISHEK	A Convolutional								
	SEHGAL AND	Neural Network								
	NASSER	Smartphone App for								$\sqrt{}$
	KEHTARNAV	Real-Time Voice								
	AZ 2021	Activity Detection								
6	Muhammad	Klasifikasi Suara Paru		./				. [
	Hasbi	Dengan		$\sqrt{}$				V		

No	Penulis / Tahun	Judul	Metode				Tujuan			
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
	Ashshiddieqy,	Convolutional Neural								
	Jondri, Achmad	Network (CNN)								
	Rizal (2020)									
7	Mazin Abed	Voice Pathology								
	Mohammed,	Detection and								
	Karrar Hameed	Classification Using								
	Abdulkareem,	Convolutional Neural								
	Salama A.	Network Model		$\sqrt{}$				$\sqrt{}$	$\sqrt{}$	
	Mostafa, Mohd									
	Khanapi Abd									
	Ghani, Mashael									
	S. Maashi,									

No	Penulis / Tahun	Judul	Metode				Tujuan			
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
	Begonya Garcia-									
	Zapirain, Ibon									
	Oleagordia,									
	Hosam									
	Alhakami and									
	Fahad Taha AL-									
	Dhief (2020)									
8	Raditya Budi	Klasifikasi Gender								
	Handoko dan	Berdasarkan Suara								
	Suyanto (2019)	Menggunakan			$\sqrt{}$		$\sqrt{}$	$\sqrt{}$		
		Support Vector								
		Machine								

No	Penulis / Tahun	Judul	Metode				Tujuan			
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
9	Aditya Singgi	Klasifikasi Suara								
	Prayogi,	Tangisan Bayi								
	Maulana Rizqi,	Berdasarkan Prosodic								
	Tresna Maulana	Features								
	Fahrudin (2019)	Menggunakan				V		V		
		Metode Moments of								
		Distribution dan K-								
		Nearest Neighbours								
10	Nur Hudha	KLASIFIKASI								
	Wijaya, Indah	SUARA JANTUNG	√`					· /		
	Soesanti, Eka	MENGGUNAKAN	V					V		
		NEURAL								

No	Penulis / Tahun	Judul			Metode	9			Tujuan	
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
	Firmansyah	NETWORK								
	(2017)	BACKPROPAGATI								
		ON BERBASIS CIRI								
		STATISTIS								
	Dede Septa	KLASIFIKASI								
	Maulana Fajar	EMOSI								
	(2023)	BERDASARKAN								
		SUARA PADA		√`			√`	$\sqrt{}$		
		DATASET CREMA-								
		D MENGGUNAKAN								
		CONVOLUTIONAL								

No	Penulis / Tahun	Judul			Metode)		Tujuan		
			ANN	CNN	SVM	KNN	MFCC	Klasifikasi	Deteksi	Aplikasi
		NEURAL								
		NETWORK								

2.2.3 Relevansi Penelitian

Relevansi Penelitian merupakan keterkaitan antar variabel pada penelitian terkait dengan penelitian yang dilakukan.

Relevansi Penelitian dapat dilihat pada Tabel 2.3.

Tabel 2.3 Relevansi Penelitian

Peneliti	(Khoirotul Aini, Budi Santoso and Dutono, 2021)	(Dede Septa Maulana Fajar, 2023)
Judul	Pemodelan CNN Untuk Deteksi Emosi Berbasis	Klasifikasi Emosi Berdasarkan Suara pada
	Speech Bahasa Indonesia	Dataset Crema-D menggunakan
		Convolutional Neural Network

Peneliti	(Khoirotul Aini, Budi Santoso and Dutono, 2021)	(Dede Septa Maulana Fajar, 2023)
Masalah	Sampai saat ini SER masih menghadapi beberapa	Penelitian dalam model klasifikasi dalam
Penelitian	tantangan seperti tingkat akurasi yang rendah dari	bidang SER telah banyak dilakukan, tetapi
	pengklasifikasi yang digunakan, kompleksitas	hanya sedikit diantara banyak penelitian
	komputasi yang tinggi, dan kelangkaan dalam	bidang SER yang menghasilkan pekerjaan
	ketersediaan natural data sets. Pengenalan emosi	yang efektif dalam memprediksi adanya emosi
	ucapan adalah tugas yang sulit karena beberapa alasan	melalui ucapan.
	seperti definisi emosi yang ambigu dan pemisahan yang	
	tidak jelas antara emosi yang berbeda.	
Objek Penelitian	Emosi berdasarkan suara manusia	Emosi berdasarkan suara manusia
Algoritma/Metod	Convolutional Neural Network	Convolutional Neural Network dan MFCC
e		
Dataset	Data set pada penelitian ini diambil dari TV series	Dataset yang digunakan dari CREMA-D yang
	berbahasa Indonesia berjudul "Imperfect". Hal ini	mencakup 7.442 klip dari 91 aktor dan aktris

Peneliti	(Khoirotul Aini, Budi Santoso and Dutono, 2021)	(Dede Septa Maulana Fajar, 2023)
	dilakukan dengan pertimbangan bahwa data dari TV	dengan usia dan etnis yang beragam dengan
	series memiliki pembagian dialog terstruktur dan	mengekspresikan enam bentuk emosi secara
	kualitas audio yang baik.	universal yaitu anger, disgust, fear, happy,
		neutral and sad.

Berdasarkan Tabel 2.3 terdapat satu penelitian yang memiliki keterkaitan dengan penelitian yang dilakukan. Penelitian yang dilakukan memiliki latar belakang masalah, objek penelitian dan metode yang digunakan cukup memiliki keterkaitan, namun dari penelitian sebelumnya yang dilakukan perlu adanya penambahan jumlah dataset yang digunakan sehingga lebih banyak data yang diklasifikasikan. Oleh karena itu, penelitian ini menggunakan dataset dari CREMA-D yang merupakan suatu standarisasi dataset untuk melakukan klasifikasi emosi berdasarkan suara.