

Mimic Human Speech in Bahasa Indonesia Using Speech Recognition and Speech Synthesis

Valens Prabagita Ivan Susilo
Department of Computer Science
President University
Cikarang, Bekasi, 17550, Indonesia
prabagita12@gmail.com

Abstrak

Sehari - hari orang menggunakan speech recognition dan speech synthesis secara tidak sadar. Teknologi ini sangat membantu dalam kegiatan mereka. Masing - masing teknologi sendiri dapat menghasilkan bermacam - macam software yang berkaitan dengan suara. Menggabungkan kedua teknologi akan menghasilkan lebih bermacam lagi software. Salah satu kombinasi itu adalah mimic speech, atau menirukan suara. Penelitian ini akan membahas tentang Speech Recognition menggunakan Convolutional Neural Network sebagai model machine learning dan Speech Synthesis yang menggunakan Concatenative Synthesis dengan suku kata sebagai satuannya. Tujuan dari penelitian ini adalah untuk mengembangkan aplikasi untuk mengumpulkan, melatih, dan menirukan suara dalam Bahasa Indonesia. Pengguna dapat berpartisipasi dengan merekam suara mereka. Suara yang dikumpulkan akan dilatih untuk digunakan dalam aplikasi untuk mengenali suara. Setelah suara yang dikumpulkan dilatih, Pengguna dapat

menggunakan aplikasi menirukan suara dengan mengidentifikasi suara dahulu lalu menghasilkan suara yang diinginkan dengan suara yang sudah teridentifikasi. Aplikasi untuk mengumpulkan dan menirukan suara akan dikembangkan dalam website dan command prompt.

1 Pendahuluan

"Ok Google, play some music". "Siri, what should I eat for lunch?". Sehari-hari orang menggunakan asisten buatan untuk membantu kegiatan mereka. Orang - orang sangat menyukai keberadaan asisten buatan ini karena dalam hitungan detik, apa yang mereka inginkan dapat terbenuhi. Secara kasat mata, orang - orang seperti berbicara dengan komputer atau ponsel mereka. Tetapi sebenarnya, speech recognition memegang peran penting didalamnya dengan bantuan machine learning. Google Assistance, Apple Siri, Microsoft Cortana, Amazon Alexa, dan asisten buatan lainnya mempunyai ribuan hingga jutaan lebih suara yang dapat dianalisis dan mereka mendapatkannya dengan mudah dari suara orang - orang yang menggunakan aplikasi mereka dengan perizinan yang mereka terima.

Jika speech recognition adalah proses untuk mendapatkan data dari menganalisa suara, kebalikannya adalah speech synthesis, yaitu proses untuk membuat suara buatan. Maka dari itu speech recognition dikenal dengan istilah speech-to-text dan speech synthesis dengan istilah text-to-speech. "Hey Cortana, read my email" adalah perintah agar asisten buatan menghasilkan suara, membacakan isi email. Masing - masing teknologi dapat menghasilkan bermacam - macam software yang berkaitan dengan suara. Menggabungkan kedua teknologi akan menghasilkan lebih bermacam lagi software. Salah satu kombinasi itu adalah mimic speech, atau menirukan suara. Penggunaan mimic speech yang dikenal adalah menciptakan suara digital yang akan digunakan sebagai suara vokal asisten buatan. Dengan itu akan membuat asisten buatan lebih personal bagi pengguna.

Penelitian ini bertujuan untuk mengembangkan aplikasi untuk mengumpulkan suara - suara, melatih machine learning dengan suara yang terkumpul, dan menirukan suara dalam Bahasa Indonesia. Aplikasi untuk mengumpulkan dan menirukan suara dikembangkan dalam website. Aplikasi dapat mengenali suara dan menghasilkan suara dari teks.

2 Batasan

Batasan aplikasi ini adalah sebagai berikut:

- Ada 9 suku kata yang dipilih untuk digunakan dalam aplikasi, a, i, na,

ma, mu, di, ri, dan ku. Suku kata akan digunakan sebagai satuan unit.

- Suara direkam dalam 1 detik dengan sample rate 16000 dan mono.

3 Metodologi

Pendekatan yang digunakan untuk mencapai tujuan penelitian ini menggunakan teknik speech synthesis concatenative synthesis, speech recognition Mel Frequency Cepstral Coefficients (MFCC), dan machine learning Convolutional Neural Network (CNN).

3.1 Concatenative Synthesis

Concatenative synthesis menghubungkan suara atau ucapan yang telah direkam. Concatenative synthesis merupakan cara termudah untuk menghasilkan suara buatan yang terdengar jelas dan alami. Salah satu aspek terpenting dalam concatenative synthesis adalah menentukan satuan suara yang benar.

Panjang dan pendek satuan mempunyai untung dan rugi masing - masing dalam penentuan ini. Dengan satuan yang panjang, lebih alami, dan titik konvergensi yang lebih sedikit, serta kontrol artikulasi yang baik terperoleh. Akan tetapi kebutuhan satuan dan memorinya akan besar. Dengan satuan yang pendek, memori yang sedikit sudah cukup, tetapi pengumpulan dan pelabelan sample suara akan lebih sulit dan kompleks [Hande, 2014].

Dalam sistem sekarang ini, satuan yang digunakan biasanya berupa kata, suku kata, *demisyllables*, fonem,

diphones, dan kadang pula triphones.

3.2 MFCC

MFCC adalah salah satu metode ekstraksi fitur yang paling umum digunakan dalam speech recognition yang diperkenalkan oleh Davis dan Mermelstein pada 1980-an [practicalcryptography.com, 2012].

3.2.1 Framing dan Windowing

Langkah yang dapat dilakukan sebelum framing dan windowing adalah menerapkan filter pre-emphasis pada sinyal untuk memperkuat frekuensi tinggi. Filter pre-emphasis dapat diterapkan ke sinyal x menggunakan filter urutan pertama dalam persamaan berikut di mana nilai tipikal untuk koefisien filter (α) adalah 0.95 atau 0.97 [haythamfayek.com, 2016]:

$$y(t) = x(t) - \alpha x(t - 1)$$

Framing dilakukan karena sinyal audio terus berubah, sehingga untuk menyederhanakan banyak hal, dengan asumsi bahwa dalam skala waktu singkat sinyal audio tidak berubah. Biasanya, framing sinyal dilakukan dalam 20-40ms bagian (25ms standar). Selang setiap bagian biasanya 10ms, yang membuat bagian satu dengan yang lainnya memiliki beberapa data yang sama. Jika, file audio tidak terbagi menjadi bagian frame yang pas, sisa bagian tersebut akan diisi dengan nol.

Setelah memotong sinyal kedalam frame, windowing atau menerapkan fungsi window seperti Hamming window kesetiap bagian frame dapat dilakukan dimana, $0 \leq n \leq N - 1$, N adalah panjang frame:

$$w[n] = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right)$$

3.2.3 Discrete Fourier Transform dan Power Spectrum

Menghitung setiap bagian dengan Discrete Fourier Transform (DFT) atau Fast Fourier Transform (FFT). Lalu dapat diikuti dengan menghitung power spectrum (periodogram). Periodogram menggunakan persamaan dimana, x_i adalah frame ke i^{th} dari sinyal x dan N adalah ukuran FFT hasil dari pangkat dua lebih besar atau sama dengan jumlah sampel dalam satu bagian frame:

$$P = \frac{|FFT(x_i)|^2}{N}$$

3.2.4 Mel Filterbank

Periodogram spektral diperkirakan masi mengandung banyak informasi yang tidak diperlukan untuk speech recognition. Mengambil segumpal tempat periodogram dan menjumlahkannya dapat mendapatkan gambaran tentang seberapa banyak energi yang ada di berbagai daerah frekuensi. Hal inilah yang dilakukan oleh mel filterbank.

Menghitung mel filterbank dengan cara menerapkan filter segitiga, biasanya sebanyak 20-40 (26 atau 40 standar) filter, pada mel-scale sampai pada power spectrum untuk mengekstrak pita frekuensi. Formula untuk mengkonversi antara Hertz (f) dan Mel (m) menggunakan persamaan berikut:

$$m = 2595 \log_{10}\left(1 + \frac{f}{700}\right)$$

$$f = 700 \left(10^{\frac{m}{2595}} - 1\right)$$

Setiap filter dalam filterbank berbentuk segitiga memiliki respons 1 pada frekuensi tengah dan menurun secara linear menuju 0 hingga mencapai frekuensi tengah dari dua filter yang berdekatan di mana responsnya adalah 0. Ukuran filter yang bagus adalah memulai filter dari 300Hz untuk yang lebih rendah dan hingga 8000Hz untuk frekuensi atas. Filterbank dapat dimodelkan dengan persamaan berikut:

$$H_m(k) = \begin{cases} 0, & k < f(m-1) \\ \frac{k-f(m-1)}{f(m)-f(m-1)}, & f(m-1) \leq k < f(m) \\ 1, & k = f(m) \\ \frac{f(m+1)-k}{f(m+1)-f(m)}, & f(m) < k \leq f(m+1) \\ 0, & k > f(m+1) \end{cases}$$

3.2.5 Logarithm

Setelah menghitung mel filterbank, selanjutnya cukup menghitung logaritmanya. Secara umum, hal ini dilakukan untuk mengandakan volume suara yang dirasakan, memasukkan 8 kali lebih banyak energi ke dalamnya.

3.2.6 Discrete Cosine Transform

Langkah terakhir adalah menghitung Discrete Cosine Transform (DCT). DCT menghias energi yang artinya matriks kovarians diagonal dapat digunakan untuk memodelkan fitur. Hanya 12 koefisien DCT yang disimpan, Ini dilakukan, karena koefisien DCT yang lebih tinggi mewakili perubahan cepat yang menurunkan kinerja speech recognition, jadi menghilangkan hingga hanya tingkat tinggi dalam neural network yang padat untuk mengenali gambar.

tersisa 12 akan mendapat peningkatan kecil.

3.3 Convolutional Neural Network

Convolutional neural network (CNN atau ConvNet) adalah salah satu varian model neural network yang dikenal untuk mengenali gambar [wikipedia.org, 2018]. Model ini dirancang untuk mengenali objek tidak peduli dimana permukaan objek itu. Model ini tidak harus mempelajari kembali perkiraan setiap permukaan yang mungkin muncul [Geitgey, 2016].

3.3.1 Convolution Layer

Convolution layer adalah lapisan untuk memasukan gambar yang telah di process sebelumnya atau output lapisan lain kedalam neural network kecil. Neural network ini memperlakukan setiap gambar atau output lapisan lain secara merata. Ia akan menandai jika sesuatu yang menarik muncul sebagai pembelajaran.

3.3.2 Max-pooling Layer

Max-pooling atau down sampling adalah lapisan untuk mengurangi output dengan cara menemukan nilai maksimum dalam output. Output dipecah menjadi kelompok dengan yang sama dan melangkah dari kelompok ke kelompok sampai mencakup seluruh output. Kemudian, setiap kelompok akan ditemukan nilai maksimum.

3.3.3 Fully-connected Layer

Fully-connected layer adalah lapisan untuk penalaran dan pertimbangan

4 Hasil Percobaan

Untuk mengevaluasi efektivitas metode yang diusulkan di bagian sebelumnya dalam aplikasi ini, percobaan dilakukan untuk memastikan aplikasi berjalan dengan baik.

Karena speech synthesis hanya menggunakan metode concatenative synthesis, atau menggabungkan, makana evaluasi akan fokus pada speech recognition. Dataset yang digunakan selama evaluasi sebanyak 10000 suara laki - laki dan perempuan, yang masing - masing 500 suara pada setiap suku kata pada suasana yang tidak berisik dan masing - masing 500 suara pada suara tidak dikenal atau suara acak. Hasil yang dianggap baik ditunjukkan dari akurasi lebih dari 75% dengan suku kata yang bersangkutan. Evaluasi juga dilakukan di suku kata diluar batasan dan suara acak seperti kondisi diam, suku kata o dan mi.

Ada 4 penguji dalam skenario ini. Tester 1 yaitu pengguna laki - laki yang suaranya telah digunakan dalam pembelajaran machine learning. Tester 2 yaitu pengguna perempuan yang suaranya telah digunakan dalam pembelajaran machine learning. Tester 3 yaitu pengguna laki - laki yang suaranya belum pernah digunakan dalam pembelajaran machine learning. Tester 4 yaitu pengguna perempuan yang suaranya belum pernah digunakan dalam pembelajaran machine learning. Setiap penguji akan melakukan evaluasi dalam keadaan seperti berikut:

- Suasana berisik (Musik keras atau orang sekitar sedang mengobrol).

- Suasana lumayan berisik (Suara hujan atau suara dari ruangan sebelah).
- Suasana tidak berisik.

Tabel 1: Hasil Tester 1 pada suasana berisik.

No	Spoken Syllable	Noisy Background			Correct Result
		1	2	3	
1	a	52 (a)	100 (a)	100 (unknown)	1/3
2	i	62 (na)	85 (ri)	96 (ma)	0/3
3	na	100 (unknown)	86 (kan)	100 (unknown)	0/3
4	ma	100 (a)	84 (unknown)	100 (a)	0/3
5	mu	99 (kan)	71 (ku)	88 (unknown)	0/3
6	di	59 (mu)	100 (ri)	97 (ri)	0/3
7	ti	50 (na)	100 (ri)	92 (a)	1/3
8	ku	100 (ku)	85 (a)	100 (kan)	1/3
9	kan	72 (unknown)	61 (kan)	79 (kan)	3/3
10	Unknown 1 (silent)	99 (unknown)	100 (a)	99 (unknown)	2/3
11	Unknown 2 (o)	100 (o)	100 (a)	100 (ku)	0/3
12	Unknown 3 (mi)	100 (a)	51 (kan)	96 (ma)	1/3

Tabel 2: Hasil Tester 1 pada suasana lumayan berisik.

No	Spoken Syllable	Semi Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	100 (a)	100 (a)	3/3
2	i	100 (i)	36 (ri)	100 (i)	2/3
3	na	97 (mu)	48 (ma)	99 (na)	1/3
4	ma	49 (i)	98 (ri)	79 (mu)	0/3
5	mu	55 (ri)	99 (ku)	50 (ri)	0/3
6	di	57 (i)	100 (di)	100 (i)	1/3
7	ti	100 (ri)	75 (ri)	100 (ri)	3/3
8	ku	92 (ku)	68 (ku)	91 (ku)	2/3
9	kan	100 (kan)	94 (na)	100 (kan)	2/3
10	Unknown 1 (silent)	100 (a)	99 (unknown)	100 (ma)	2/3
11	Unknown 2 (o)	100 (ku)	89 (ku)	72 (ku)	1/3
12	Unknown 3 (mi)	100 (ri)	94 (ri)	100 (ri)	0/3

5 MIMIC HUMAN SPEECH

Tabel 3: Hasil Tester 1 pada suasana tidak berisik.

No	Spoken Syllable	Not Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	100 (a)	100 (a)	3/3
2	i	94 (i)	48 (kan)	94 (i)	2/3
3	na	80 (na)	99 (na)	100 (na)	3/3
4	ma	96 (di)	90 (i)	57 (di)	0/3
5	mu	100 (i)	86 (ku)	95 (i)	0/3
6	di	99 (ri)	100 (i)	80 (ri)	0/3
7	ti	100 (ri)	96 (ri)	91 (ri)	3/3
8	ku	85 (ku)	96 (ku)	100 (ku)	3/3
9	kan	92 (kan)	100 (kan)	100 (kan)	3/3
10	Unknown 1 (silent)	98 (unknown)	98 (unknown)	98 (unknown)	3/3
11	Unknown 2 (o)	100 (o)	98 (ku)	65 (ku)	1/3
12	Unknown 3 (mi)	54 (mu)	100 (ri)	67 (mu)	2/3

Tabel 6: Hasil Tester 2 pada suasana tidak berisik.

No	Spoken Syllable	Not Noisy Background			Correct Result
		1	2	3	
1	a	88 (a)	99 (a)	99 (a)	3/3
2	i	99 (i)	99 (i)	52 (i)	2/3
3	na	97 (a)	54 (kan)	99 (a)	0/3
4	ma	99 (na)	99 (na)	98 (na)	0/3
5	mu	81 (i)	95 (i)	44 (i)	0/3
6	di	99 (i)	36 (i)	96 (i)	0/3
7	ti	47 (ri)	32 (unknown)	48 (i)	0/3
8	ku	98 (ku)	67 (kan)	99 (ku)	2/3
9	kan	84 (na)	85 (na)	87 (na)	0/3
10	Unknown 1 (silent)	99 (unknown)	99 (unknown)	100 (unknown)	3/3
11	Unknown 2 (o)	100 (o)	99 (ku)	95 (ku)	0/3
12	Unknown 3 (mi)	67 (ri)	78 (di)	77 (i)	1/3

Tabel 4: Hasil Tester 2 pada suasana berisik.

No	Spoken Syllable	Noisy Background			Correct Result
		1	2	3	
1	a	99 (unknown)	99 (a)	99 (a)	2/3
2	i	82 (unknown)	98 (ma)	99 (ri)	0/3
3	na	52 (ri)	99 (ma)	61 (ma)	0/3
4	ma	50 (unknown)	37 (unknown)	95 (ri)	0/3
5	mu	99 (unknown)	52 (di)	99 (unknown)	0/3
6	di	99 (i)	99 (unknown)	99 (unknown)	0/3
7	ti	86 (i)	20 (kan)	57 (di)	0/3
8	ku	99 (ku)	68 (kan)	90 (unknown)	1/3
9	kan	84 (kan)	99 (kan)	78 (a)	2/3
10	Unknown 1 (silent)	100 (unknown)	100 (unknown)	100 (unknown)	3/3
11	Unknown 2 (o)	99 (ma)	99 (ku)	99 (ku)	0/3
12	Unknown 3 (mi)	36 (di)	96 (mu)	99 (di)	1/3

Tabel 7: Hasil Tester 3 pada suasana berisik.

No	Spoken Syllable	Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	99 (a)	100 (a)	3/3
2	i	71 (ma)	99 (ri)	99 (ri)	0/3
3	na	70 (ku)	99 (a)	57 (unknown)	0/3
4	ma	99 (a)	100 (a)	100 (a)	0/3
5	mu	74 (mu)	96 (kan)	74 (mu)	0/3
6	di	99 (ri)	96 (ri)	89 (ri)	0/3
7	ti	100 (unknown)	90 (na)	86 (a)	0/3
8	ku	98 (kan)	64 (mu)	99 (a)	0/3
9	kan	99 (na)	100 (kan)	98 (kan)	2/3
10	Unknown 1 (silent)	99 (a)	99 (a)	100 (a)	0/3
11	Unknown 2 (o)	92 (na)	84 (na)	61 (mu)	1/3
12	Unknown 3 (mi)	64 (kan)	53 (na)	99 (ri)	2/3

Tabel 5: Hasil Tester 2 pada suasana lumayan berisik.

No	Spoken Syllable	Semi Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	100 (unknown)	99 (a)	2/3
2	i	92 (ri)	97 (unknown)	96 (ri)	0/3
3	na	83 (ma)	62 (unknown)	75 (kan)	0/3
4	ma	87 (unknown)	64 (ma)	90 (na)	2/3
5	mu	76 (ma)	78 (ku)	66 (di)	0/3
6	di	58 (mu)	48 (mu)	30 (kan)	0/3
7	ti	99 (kan)	89 (kan)	99 (ri)	1/3
8	ku	99 (unknown)	99 (ku)	88 (unknown)	1/3
9	kan	99 (kan)	65 (kan)	85 (kan)	2/3
10	Unknown 1 (silent)	100 (unknown)	79 (a)	99 (ri)	1/3
11	Unknown 2 (o)	99 (a)	82 (ku)	99 (na)	0/3
12	Unknown 3 (mi)	89 (kan)	35 (ku)	96 (ri)	1/3

Tabel 8: Hasil Tester 3 pada suasana lumayan berisik.

No	Spoken Syllable	Semi Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	99 (a)	99 (a)	3/3
2	i	99 (ri)	90 (di)	92 (di)	0/3
3	na	100 (kan)	99 (kan)	92 (ma)	0/3
4	ma	99 (a)	100 (kan)	100 (kan)	0/3
5	mu	99 (mu)	99 (mu)	65 (ku)	2/3
6	di	99 (ri)	100 (ri)	99 (ri)	0/3
7	ti	95 (kan)	76 (i)	78 (di)	0/3
8	ku	99 (a)	96 (a)	100 (a)	0/3
9	kan	99 (a)	96 (a)	99 (a)	0/3
10	Unknown 1 (silent)	99 (ri)	73 (di)	60 (unknown)	2/3
11	Unknown 2 (o)	100 (a)	99 (ku)	100 (a)	0/3
12	Unknown 3 (mi)	99 (na)	99 (ri)	99 (i)	0/3

Tabel 9: Hasil Tester 3 pada suasana tidak berisik.

No	Spoken Syllable	Not Noisy Background			Correct Result
		1	2	3	
1	a	100 (a)	100 (a)	100 (a)	3/3
2	i	92 (ri)	95 (ri)	41 (ku)	0/3
3	na	99 (na)	74 (ku)	80 (a)	1/3
4	ma	59 (a)	84 (ku)	48 (kan)	0/3
5	mu	99 (mu)	99 (ri)	99 (ri)	1/3
6	di	100 (ri)	99 (i)	99 (i)	0/3
7	ti	99 (ri)	96 (ri)	94 (i)	2/3
8	ku	73 (di)	100 (mu)	99 (mu)	0/3
9	kan	100 (kan)	99 (kan)	96 (kan)	3/3
10	Unknown 1 (silent)	98 (unknown)	98 (unknown)	98 (unknown)	3/3
11	Unknown 2 (o)	99 (ku)	68 (mu)	100 (ku)	1/3
12	Unknown 3 (mi)	99 (ri)	93 (i)	100 (ri)	0/3

Tabel 10: Hasil Tester 4 pada suasana berisik.

No	Spoken Syllable	Noisy Background			Correct Result
		1	2	3	
1	a	100 (unknown)	99 (unknown)	94 (ma)	0/3
2	i	97 (unknown)	99 (ri)	99 (ri)	0/3
3	na	100 (na)	99 (ri)	99 (ma)	1/3
4	ma	93 (na)	99 (ma)	100 (na)	1/3
5	mu	55 (ri)	99 (unknown)	43 (unknown)	0/3
6	di	99 (ri)	54 (i)	99 (ma)	0/3
7	ti	97 (ri)	74 (ri)	100 (ri)	2/3
8	ku	68 (ma)	92 (ku)	52 (s2)	1/3
9	kan	48 (ku)	22 (ku)	99 (a)	0/3
10	Unknown 1 (silent)	100 (unknown)	100 (unknown)	100 (unknown)	3/3
11	Unknown 2 (o)	100 (a)	100 (a)	100 (a)	0/3
12	Unknown 3 (mi)	55 (mu)	99 (ri)	94 (i)	1/3

Tabel 11: Hasil Tester 4 pada suasana lumayan berisik.

No	Spoken Syllable	Semi Noisy Background			Correct Result
		1	2	3	
1	a	73 (a)	91 (unknown)	99 (unknown)	0/3
2	i	81 (ri)	99 (i)	100 (ri)	1/3
3	na	100 (ma)	99 (unknown)	99 (unknown)	0/3
4	ma	99 (ma)	99 (na)	78 (ma)	2/3
5	mu	81 (unknown)	99 (unknown)	99 (unknown)	0/3
6	di	100 (ri)	99 (i)	99 (i)	0/3
7	ti	85 (ri)	99 (ri)	98 (ri)	3/3
8	ku	96 (ku)	89 (ku)	99 (a)	2/3
9	kan	100 (a)	99 (a)	64 (ma)	0/3
10	Unknown 1 (silent)	99 (unknown)	99 (ri)	100 (unknown)	2/3
11	Unknown 2 (o)	99 (a)	100 (a)	99 (a)	0/3
12	Unknown 3 (mi)	43 (a)	59 (ri)	99 (unknown)	3/3

Tabel 12: Hasil Tester 4 pada suasana tidak berisik.

No	Spoken Syllable	Not Noisy Background			Correct Result
		1	2	3	
1	a	99 (unknown)	75 (a)	100 (a)	2/3
2	i	88 (i)	99 (i)	99 (ri)	2/3
3	na	99 (na)	99 (ma)	99 (ma)	1/3
4	ma	55 (ma)	72 (kan)	99 (na)	0/3
5	mu	99 (unknown)	94 (unknown)	28 (i)	0/3
6	di	89 (i)	99 (ri)	99 (ri)	0/3
7	ti	99 (ri)	100 (ri)	99 (ri)	3/3
8	ku	39 (kan)	99 (ku)	58 (mu)	1/3
9	kan	100 (a)	99 (a)	100 (a)	3/3
10	Unknown 1 (silent)	40 (unknown)	24 (unknown)	60 (unknown)	3/3
11	Unknown 2 (o)	100 (a)	99 (a)	100 (a)	0/3
12	Unknown 3 (mi)	60 (ri)	100 (i)	99 (ri)	1/3

Meskipun suasana berisik dan lumayan berisik, setiap tester tidak menunjukkan hasil yang cukup baik. Masih dapat memprediksi sekitar 2 atau 3 suku kata dengan benar dan akurasi yang tinggi. Suasana tidak berisik juga tidak menjamin hasil yang sempurna bahkan pada penguji 1 dan 2. Tetapi, hasilnya lebih baik daripada suasana berisik dan lumayan berisik. Selain itu, pengurangan suasana latar belakang tidak diterapkan ketika proses ekstraksi yang dilakukan sebelum pembelajaran.

Penguji 1 dan 2 seringkali mendapat hasil yang lebih baik dari pada penguji 3 dan 4. Tetapi, pada saat tertentu penguji 3 dan 4 mendapat hasil yang lebih baik. Ini bisa terjadi dikarenakan suasana latar belakangnya atau pun juga kemampuan mikrofon untuk merekam suara.

Suku kata ma, mu, dan di, serta suku kata acak o dan mi memiliki hasil yang buruk hampir setiap waktu. Suku kata acak o dan mi dikarekanan pada pembelajarannya suara acak berisikan kebanyakan suasana latar belakang dari pada suku kata o dan mi. Model machine learning tidak dapat

memprediksi suku kata *ma*, *mu*, dan *di* dengan baik. Seringkali model memprediksi dengan benar tapi dengan akurasi yang rendah atau memprediksi suku kata relatif, contoh hasil *ri* atau *i* pada suku kata *di*. Hasil ini dapat ditingkatkan dengan menambahkan data lagi untuk dipelajari karena data yang digunakan terhitung masih relatif sedikit. Model machine learning yang optimum juga dapat meningkatkan hasil prediksi dan akurasi.

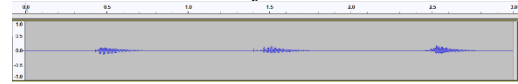
Berikut adalah teks acak untuk mengetes dan mengevaluasi suara yang dibuat atau dihasilkan:

- *diriku*
- *aku makan ikan*
- *di mana mamamu*
- *halo namaku ivan*

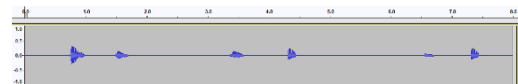
'*diriku*', '*aku makan ikan*', dan '*di mana mamamu*' dapat dihasilkan. Namun, '*halo namaku ivan*' tidak bisa. Ini dikarenakan '*h*' tidak ditemukan dalam database saat teks sedang dianalisa dari depan. Aplikasi akan menampilkan peringatan di browser bahwa '*h*' tidak ditemukan dan akan menampilkan daftar suku kata yang sudah teridentifikasi. Sebenarnya '*h*' itu sendiri tidak termasuk dalam suku kata pilihan. Teks yang berisikan suku kata yang tidak terpilih maupun terpilih namun belum dianalisa dan didaftarkan maka aplikasi akan menampilkan peringatan. Suku kata apapun yang tidak ditemukan dalam database akan menghentikan proses membuat suara.

Berikut adalah hasil pembuatan suara dari '*diriku*', '*aku makan ikan*', dan '*di mana mamamu*' secara berurutan:

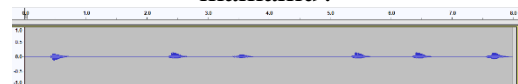
Gambar 1: *Waveform* dari '*diriku*'.



Gambar 2: *Waveform* dari '*aku makan ikan*'.

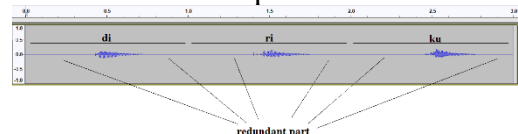


Gambar 3: *Waveform* dari '*di mana mamamu*'.



Seperti yang ditunjukkan pada gambar – gambar diatas, setiap teks dapat dibedah dan dibedakan secara mudah. Setiap suara direkam dalam 1 detik, maka durasi suara yang dihasilkan dari masing – masing teks adalah jumlah dari suku kata pada teks tersebut. Meskipun suara yang dihasilkan baik, dapat didengar dan dimengerti, namun masih ada bagian hening atau suara yang tidak diperlukan kecuali spasi. Hal itu membuat suara yang dihasilkan kurang terdengar lancer. Hal ini dapat terjadi dikarenakan tidak ada proses lebih lanjut untuk menganalisa dan menghapus bagian suara yang tidak diperlukan pada proses penggabungan ataupun setelah proses identifikasi dalam aplikasi ini.

Figure 4: *Waveform* dari '*diriku*' menunjukan suku kata dan suara yang tidak diperlukan.



5 Diskusi

Pada bagian ini, ada beberapa diskusi mengapa metode yang diusulkan digunakan untuk penelitian ini. Poin - poin penting adalah:

- Concatenative Synthesis untuk Speech Synthesis
- Suku kata untuk satuan suara
- MFCC untuk Speech Recognition
- CNN untuk Machine Learning model

Dalam mimic speech, suara diambil dari suara yang telah dianalisa. Concatenative synthesis dapat menjadi pendekatan terbaik untuk masalah ini. Selain caranya paling mudah untuk digunakan, cara ini juga cepat untuk dikembangkan kedalam aplikasi. Articulatory dan Formant synthesis terlalu kompleks karena parameter yang dibutuhkan banyak untuk mengembangkan system vokal suara ataupun peraturan yang dapat menyesuaikan semua pengguna.

Sangat sulit untuk menemukan penelitian mengenai jumlah pasti *demisyllables*, fonem atau satuan suara yang lebih rendah dalam Bahasa Indonesia. Maka dari itu suku kata adalah pilihan yang paling tepat karena kata akan terlalu banyak memakan memori dan kurang fleksibel dalam membuahkan suara dalam bentuk kalimat.

MFCC mempunyai keuntungan dibandingkan dengan Linear Prediction Coefficients (LPC) [Dave & Pipalia, 2014] yaitu karena MFCC dapat menirukan sistem pendengaran manusia dengan baik. Meskipun Perceptually Based Linear Predictive

Analysis (PLP) juga dapat melakukan hal yang sama, MFCC tetap akan dipilih karena MFCC adalah proses ekstraksi yang paling sering digunakan. Maka dari itu, teknik tersebut sudah banyak tersebar menjadikan mudah untuk dikembangkan dan *debug* dalam aplikasi.

Dua alasan utama menggunakan CNN. Pertama, concatenative synthesis menggunakan satuan suara. Tiap satuan suara adalah bagian kecil dari suatu kalimat atau kata. CNN dianggap sebagai pendekatan optimal untuk menganalisa *small-footprint keyword*, yakni satuan suara itu sendiri, dari pada pendekatan machine learning yang lain [Sainath & Parada, 2015]. Kedua, proses ekstraksi MFCC dapat di plot menjadi spektogram. Dan CNN adalah pendekatan paling umum untuk menyelesaikan masalah tentang visual dan gambar, termasuk spektogram.

6 Conclusion

Ada beberapa kesimpulan yang dapat diperoleh dari penelitian ini:

- Aplikasi ini dapat mengumpulkan data suara melalui website.
- Aplikasi ini dapat melatih model machine learning dengan data yang sudah dikumpulkan melalui command prompt.
- Aplikasi ini dapat menirukan suara Bahasa Indonesia melalui website.
- Aplikasi ini dapat mengenali suara dari rekaman suara walaupun prediksi dan akurasi tidak sempurna. Akan tetapi, machine learning dapat memprediksi hampir setiap waktu dengan baik.

- Aplikasi ini dapat menghasilkan suara berdasarkan text yang dimasukan walaupun hasilnya masih memiliki bagian yang hening atau suara yang tidak diperlukan. Akan tetapi, hasil yang dihasilkan dapat didengar dan dimengerti.

Di masa yang akan datang, penelitian lebih lanjut dalam speech recognition pada bagian meningkatkan model machine learning. Tidak ada benar dalam memodelkan machine learning, tetapi selalu ada model yang optimum untuk mendapatkan prediksi dan akurasi yang paling baik. Dalam speech synthesis dengan menambahkan proses untuk menghilangkan bagian hening atau suara tidak perlu dari rekaman suara dan juga menghilangkan atau mengurangi suasana latar belakang akan menghasilkan pembuatan suara yang lebih lancar dan baik didengar. alam Bahasa Indonesia, penelitian untuk menentukan fonem Bahasa Indonesia dapat menjadi peningkatan yang signifikan dalam speech recognition dan speech synthesis karena aplikasi ini menggunakan suku kata sebagai satuan suara.

7 Ucapan Terima Kasih

Penulis hendak mengucapkan ucapan terima kasih kepada Bapak Tjong Wan Sen sebagai penasihat skripsi untuk saran dan dukungan selama proses skripsi ini dilaksanakan. Dan juga kepada dosen - dosen yang lain untuk dukungan mereka, ilmu, dan pengalaman selama kuliah.

Referensi

- [Hande, 2014] Hande, S. S. (2014). A Review on Speech Synthesis an Artificial Voice Production. *International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering*, 8.
- [practicalcryptography.com, 2012] practicalcryptography.com. (2012). *Mel Frequency Cepstral Coefficient (MFCC) tutorial*. Retrieved from Practical Cryptography: <http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/>
- [haythamfayek.com, 2016] haythamfayek.com. (2016, April 21). *Speech Processing for Machine Learning: Filter banks, Mel-Frequency Cepstral Coefficients (MFCCs) and What's In-Between*. Retrieved from Speech Processing for Machine Learning: Filter banks, Mel-Frequency Cepstral Coefficients (MFCCs) and What's In-Between: <https://haythamfayek.com/2016/04/21/speech-processing-for-machine-learning.html>
- [wikipedia.org, 2018] wikipedia.org. (2018, December 26). *Convolutional neural network*. Retrieved from Wikipedia: https://en.wikipedia.org/wiki/Convolutional_neural_network
- [Geitgey, 2016] Geitgey, A. (2016, Juny 14). *Machine Learning is Fun! Part 3: Deep Learning and Convolutional Neural Networks*. Retrieved from Medium:

<https://medium.com/@ageitgey/machine-learning-is-fun-part-3-deep-learning-and-convolutional-neural-networks-f40359318721>

[Sainath & Parada, 2015] Sainath, T. N., & Parada, C. (2015). *Convolutional Neural Networks for Small-footprint Keyword Spotting*. New York, NY, U.S.A: Google, Inc.