



### **Hak cipta dan penggunaan kembali:**

Lisensi ini mengizinkan setiap orang untuk menggubah, memperbaiki, dan membuat ciptaan turunan bukan untuk kepentingan komersial, selama anda mencantumkan nama penulis dan melisensikan ciptaan turunan dengan syarat yang serupa dengan ciptaan asli.

### **Copyright and reuse:**

This license lets you remix, tweak, and build upon work non-commercially, as long as you credit the origin creator and license it on your new creations under the identical terms.

## BAB II

### TINJAUAN PUSTAKA

#### 2.1. Landasan Teori

Pada penelitian ini, sebuah aplikasi *speech recognition* kamus geografi dirancang menggunakan metode *Convolutional Neural Network* berbasis aplikasi komputer. Penelitian terdahulu adalah penelitian oleh Edbert Hansel [2] dengan membuat aplikasi *speech recognition* untuk kamus geografi berbahasa Indonesia. Abdel-Hamid [1] berhasil membuktikan bahwa CNN mempunyai yang sangat baik dalam mengenali suara yang dapat diterapkan ke aplikasi *speech recognition*.

Analisa metode CNN ini dibuat pada sistem berbasis aplikasi *desktop*. Input adalah hasil rekaman yang diucapkan oleh user. Proses pelatihan kata menggunakan metode *Convolutional Neural Network* (CNN) yang diterapkan pada framework TensorFlow. Rancangan model CNN dalam aplikasi dibuat berdasarkan penelitian yang dilakukan oleh Tara N. Sainath dan Carolina Parada [3] dari Google pada konferensi INTERSPEECH 2015. Kelebihan aplikasi ini adalah lebih akurat dalam fitur *speech recognition*, sehingga hasil yang dikeluarkan akurat dan anak-anak Indonesia lebih tertarik untuk mempelajari Geografi.

#### 2.2. TensorFlow

TensorFlow merupakan *library open source* untuk perhitungan numerik menggunakan *data flow graphs*. *Nodes* pada graph mewakili operasi matematika, sedangkan *graph edges* mewakili array multidimensi (tensor) yang mengalir dari keduanya. TensorFlow awalnya dikembangkan oleh *researchers* dan *engineers* dari

Google Brain di dalam Machine Intelligence Research Google, dengan tujuan melakukan riset *machine learning* dan *deep neural networks* [4].

TensorFlow mendukung pengembangan aplikasi dengan frameworknya di sistem operasi Windows dan Linux. Semua proses di dalam TensorFlow melibatkan tensor. Tensor adalah generalisasi dari vector atau matriks untuk mencapai dimensi yang lebih tinggi. Secara internal, TensorFlow membuat tensor sebagai array n-dimensi yang berisi tipe data dasar. Semua elemen di dalam Tensor mempunyai tipe data yang sama dan tipe datanya harus diketahui [4].

TensorFlow dapat digunakan untuk melakukan algoritma *machine learning* dan dapat mengimplementasikan algoritma tersebut [5]. Perhitungan di dalam TensorFlow dapat dieksekusi dalam sistem yang sangat luas, dari perangkat *mobile* seperti handphone hingga komputer canggih dengan perhitungan kartu GPU. Karena sistem yang fleksibel, banyak algoritma yang rumit dapat ditulis ke dalam TensorFlow, termasuk pelatihan dari algoritma *Deep Neural Network*.

### 2.2.1. Tensor

Tensor adalah bagian inti dari TensorFlow. Tensor dapat mempunyai *rank* untuk tiap dimensi yang dipunyai. Biasanya disebut *order* atau *degree*. Rank tersebut dapat dilihat pada tabel 2.1.

Tabel 2.1 Rank TensorFlow [4]

Rank	Entitas Matematika
0	Scalar (hanya besaran)
1	Vector (besaran dan arah)
2	Matriks (tabel dari angka)
3	3-Tensor (kubus dari angka)
N	n-Tensor

Tensor dapat diisi berbagai macam tipe data. *Signed* atau *unsigned* integer dari 8 bits sampai 64 bits. Float dan double IEEE, angka kompleks, string juga dapat dimasukkan.

### 2.2.2. Variable

Variable di dalam TensorFlow adalah tempat untuk menyimpan *persistent and shared state* variabel yang nantinya dimanipulasi oleh program. Cara terbaik membuat variabel adalah dengan menggunakan fungsi `tf.get_variable()`. Fungsi ini digunakan untuk mengakses sebuah variable dari banyak replika. Selain itu dengan variabel, variabel yang sudah pernah dibuat untuk dapat digunakan kembali [4].

### 2.2.3. Operasi dan Kernel

Operasi di dalam TensorFlow mempunyai nama dan mewakili komputasi yang abstrak seperti penambahan atau perkalian matriks. Operasi dapat mempunyai atribut dan semua atribut harus diberikan saat pembangunan *graph* agar dapat membuat *node* melakukan sebuah operasi [5]. Kegunaan atribut yang paling umum

adalah untuk membuat operasi antar 2 tipe variabel yang berbeda. Contohnya adalah penambahan 2 tensor bertipe float ke dalam 2 tensor bertipe int32.

Kernel adalah implementasi operasi di dalam *device* tertentu. Bisa dijalankan dalam CPU ataupun GPU. TensorFlow *binary* dapat mendefinisikan kumpulan operasi dan kernel yang tersedia melalui mekanisme registrasi. Sehingga jika diperlukan operasi atau kernel dari luar, dapat ditambahkan via registrasi.

#### 2.2.4. Sessions

Klien berinteraksi dengan sistem TensorFlow melalui pembuatan *Session*. Untuk membuat *computation graph*, antar muka *Session* mendukung *method Extend* untuk menggabungkan *Session* yang sedang dijalankan dengan tambahan *nodes* atau *edges*. Operasi utama lain yang didukung oleh *Session* adalah *Run*, dengan mengambil kumpulan nama keluaran yang harus dihitung dan kumpulan tensor yang harus diberikan ke dalam *graph* untuk menghasilkan beberapa keluaran *nodes*. Dengan *Run*, TensorFlow dapat menjalankan semua perhitungan *nodes* dengan tujuan mendapatkan keluaran yang diminta [5].

### 2.3. Speech Recognition

*Automatic Speech Recognition* oleh mesin sudah menjadi tujuan riset lebih dari empat dekade dan telah menjadi inspirasi untuk banyak film *sci-fi* seperti komputer HAL di film “A Space Odyssey” dan robot R2D2 di film “Star Wars”. Tetapi untuk mencapai tujuan tersebut, para ilmuwan masih jauh dari memperoleh mesin yang dapat mengenali setiap subjek dari semua orang di dalam segala kondisi [6]. Teknologi ini memungkinkan sebuah mesin untuk dapat mengolah suara dari

seseorang dan memahami kata yang diucapkan. Cara kerjanya adalah dengan mengubah suara orang menjadi sinyal digital yang berisi angka, kemudian angka tersebut disesuaikan dengan pola yang sudah dipelajari untuk menentukan kata yang diucapkan oleh seseorang. Nantinya mesin juga dapat menampilkan keluaran berupa bentuk teks ke user.

Salah satu aspek tersulit untuk melakukan riset dalam *speech recognition* oleh mesin adalah karena diperlukan lebih dari satu cabang ilmu dan langkah yang biasanya dilakukan oleh periset adalah menerapkan pendekatan monolitik ke masalah individu. Beberapa disiplin ilmu yang sudah pernah diterapkan ke permasalahan *speech recognition* [6]:

- **Pengolahan Sinyal:** Pengolahan Sinyal adalah proses ekstraksi informasi yang relevan dari sinyal suara. Di dalam pengolahan sinyal, *spectral analysis* termasuk salah satu contoh penerapan untuk mengetahui properti *time-varying* dari sinyal suara. Beberapa tipe *preprocessing* sinyal juga termasuk dalam membuat sinyal suara bagus untuk diolah nantinya.
- **Physics (acoustics):** *Physics* adalah ilmu memahami hubungan antara sinyal suara fisik dan mekanisme fisiologi (mekanisme sistem vokal manusia) yang menghasilkan suara yang dapat didengar (mekanisme sistem pendengaran manusia).
- **Pattern Recognition:** *Pattern Recognition* adalah kumpulan algoritma yang digunakan untuk mengelompokkan data sesuai

dengan polanya (*clustering*) dan mencocokkan pasangan pola untuk dibandingkan fitur polanya.

- Teori Komunikasi dan Informasi: Teori ini adalah prosedur untuk mengestimasi parameter dari model statistik. Metode untuk mendeteksi keberadaan pola suara tertentu. Ini juga merupakan algoritma *coding* dan *decoding*. Digunakan untuk mencari jalan terbaik dari *grid* yang jumlahnya terbatas untuk menentukan urutan kata terbaik yang dikenali.
- *Linguistics*: *Linguistics* adalah hubungan antar suara (*phonology*), kata di dalam bahasa (*syntax*), arti dari kata yang diucapkan (*semantics*), dan *sense* yang didapatkan dari arti (*pragmatics*). Di dalam ilmu ini adalah metodologi dari *grammar* dan pengolahan kata.
- Fisiologi: Fisiologi adalah pemahaman dari mekanisme tingkat tinggi dari sistem saraf manusia untuk menghasilkan suara dan menghasilkan persepsi di manusia. Banyak teknik modern yang mencoba menerapkan pengetahuan ini di *framework artificial neural networks*.
- Ilmu Komputer: Dampak dari Ilmu Komputer adalah terhadap pembelajaran algoritma yang efisien untuk diimplementasikan, di *software* maupun *hardware*, juga bermacam metode yang digunakan untuk praktek sistem pengenalan suara.



- Psikologi: Ilmu yang mempelajari faktor yang memungkinkan teknologi untuk digunakan oleh manusia dalam kehidupan sehari-hari.

Dalam pendekatan *Automatic Speech Recognition* oleh mesin, ada 3 jenis pendekatan[6]:

- Pendekatan *acoustic-phonetic*
- Pendekatan *pattern recognition*
- Pendekatan *artificial intelligence*

Pendekatan melalui *acoustic-phonetic* sudah dipelajari lama dan lebih dari 40 tahun. Tetapi karena beberapa alasan, pendekatan ini belum mencapai kesuksesan dibandingkan pendekatan alternatif lainnya. Pendekatan ini berpusat pada teori akustik fonetik yang membedakan tiap unit fonetik dalam sebuah bahasa dan unit fonetik ini mempunyai karakteristik dalam sinyal suaranya, atau dalam spectrumnya. Fase pertamanya adalah segmentasi dan *labeling* dimana fase ini memecah sinyal suara ke bentuk *discrete*, dimana properti akustik dari sebuah sinyal mewakili sebuah unit fonetik, dan memberi label fonetik ke bagian yang dipecah sesuai properti akustiknya. Fase kedua adalah menebak kata dari kumpulan unit fonetik yang didapat di fase satu[6].

Pendekatan melalui *pattern-recognition* mempunyai dua langkah, yaitu pelatihan dari pola suara dan memahami pola melalui perbandingan pola. Pemahaman akan suara diberikan ke sistem melalui prosedur pelatihan. Algoritma dibuat dengan *training set* yang disediakan. Pelatihan ini seharusnya bisa



mengenali properti akustik dari sebuah pola. Mesin dapat melakukan *pattern classification* karena mesin mempelajari properti akustik mana yang dapat diandalkan selama dalam proses pelatihan. Pendekatan ini adalah pendekatan yang cocok untuk pengenalan suara karena tiga alasan[6]:

- Pemakaian yang mudah
- Ketahanan dan ketetapan untuk kamus, pengguna, set fitur, algoritma pembandingan, dan aturan pengambilan keputusan yang berbeda.
- Terbukti sangat bagus

Sedangkan pendekatan *artificial intelligence* adalah pendekatan yang menggabungkan pendekatan *acoustic-phonetic* dan pendekatan *pattern recognition*. Pendekatan ini mengeksplorasi konsep kedua pendekatan lainnya. Pendekatan ini mencoba menjalankan prosedur pengenalan seperti seseorang menerapkan kepiatarannya dalam memvisualisasi, menganalisa, dan akhirnya membuat keputusan berdasarkan fitur akustik. Sehingga teknik yang digunakan dalam metode ini meliputi penggunaan sistem ahli dalam segmentasi dan labeling, sehingga langkah paling krusial dan tersulit dapat dilakukan tidak hanya dengan informasi akustik di pendekatan *acoustic-phonetic* saja. Tetapi pendekatan ini juga dapat belajar dan beradaptasi seiring waktu dan menggunakan *neural networks* untuk mempelajari hubungan antar fonetik, *input*, dan diskrimasi antar kelas suara yang mirip[6].

Penggunaan *neural networks* dapat dikatakan sebagai implementasi arsitektur yang dapat diterapkan dalam tiga pendekatan klasik. Konsep dan ide menerapkan *neural networks* ke permasalahan pengenalan suara masih baru-baru ini[6]. Sehingga *neural networks* dianggap sebagai kunci menyelesaikan masalah pengenalan suara di beberapa kasus nyata. Penelitian ini nantinya juga berpusat pada penggunaan *neural networks* sebagai dasar arsitektur.

#### 2.4. Convolutional Neural Network

Convolutional Neural Network (CNN) menyediakan *shift invariance* terhadap waktu dan cakupan dan penting untuk mencapai performa *state-of-the-art* untuk image recognition[7]. Tetapi CNN juga sudah terbukti dapat meningkatkan akurasi pengenalan suara disbanding *Deep Neural Network* murni di beberapa tugas tertentu [1], [8]. Untuk mendukung CNN, beberapa titik komputasi perlu diimplementasikan.

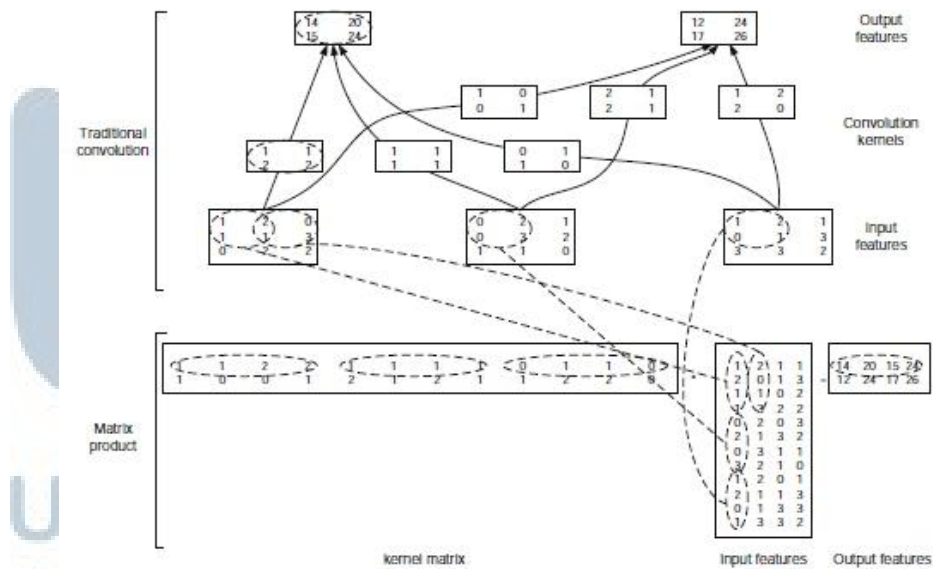
Di penelitian sebelumnya penggabungan *deep neural network* – *hidden Markov Model* (DNN-HMM) menunjukkan perkembangan pada performa pengenalan suara dibandingkan model konvensional *Gaussian mixture model* (GMM-HMM). Penambahan pada performa adalah karena kapasitas DNN untuk memodelkan korelasi kompleks pada fitur suara. Terlebih lagi, CNN dapat menurunkan *error rate reduction*[1].

*Convolution* dalam kasus ini adalah proses konvolusi produk berelemen dari sebuah kernel menjadi sebuah gambar. *Input* dari *convolution node* mempunyai 3 kanal (direpresentasikan sebagai matriks 3x3) dan keluarannya mempunyai 2 kanal

(direpresentasikan sebagai matriks 2x2). Kanal adalah pandangan dari gambar yang sama. Contohnya sebuah gambar RGB bisa direpresentasikan menjadi tiga kanal: R, G, dan B, setiap kanal ukurannya sama[7].

Dalam setiap pasangan kanal *input* dan *output* ada kernelnya. Jumlah kernel sama dengan jumlah produk dari kanal *input*  $C_x$  dan jumlah kanal *output*  $C_y$ . Pada Gambar 2.1 di bawah,  $C_x = 3$ ,  $C_y = 2$  dan jumlah kernelnya adalah 6. Setiap kernel  $K_{kl}$  dari kanal *input*  $k$  dan kanal *output*  $l$  adalah matriks. Kernel bergerak mengikuti *input* dengan *subsampling rate*  $S_r$  dan  $S_c$  berurutan pada arah vertikal dan horizontal. Untuk setiap *output* kanal  $l$  dan irisan *input*  $(i,j)$ , dimana  $Y_{kij}$  mempunyai ukuran yang sama dengan  $K_{kl}$ .

$$v_{lij}(K, Y) = \sum_k vec(K_{kl}) \circ vec(Y_{kij}) \quad (2.1)$$



Gambar 2.1 Contoh Operasi Konvolusi[7]

Fungsi evaluasi ini melibatkan banyak operasi matriks kecil dan bisa lambat. Chellapila et al. [9] mengajukan teknik untuk mengkonversi semua operasi matriks kecil menjadi sebuah matriks produk besar yang terlihat di bawah Gambar 2.1. Dengan ini, semua parameter kernel dapat digabungkan menjadi sebuah kernel matriks besar  $\mathbf{W}$  seperti yang terlihat di bawah kiri Gambar 2.1. Agar titik konvolusi dapat digunakan oleh titik konvolusi yang lain, konversi harus dibuat sedikit berbeda dari yang dirancang oleh Chellapila et al. [9] dengan melakukan *transpose* kedua kernel matriks dan fitur input matriks serta orde matriks di produk. Dengan perubahan ini, setiap sampel di keluaran dapat direpresentasikan dengan kolom  $O_r \times O_c$  dari  $C_v \times 1$  vector yang dapat dirubah bentuknya menjadi sebuah kolom, dimana

$$O_r = \begin{cases} \frac{I_r - K_r}{S_r} + 1 & \text{no padding} \\ \left( \frac{I_r - \text{mod}(K_r, 2)}{S_r} \right) + 1 & \text{zero padding} \end{cases} \quad (2.2)$$

adalah jumlah baris dari *output* gambar, dan

$$O_c = \begin{cases} \frac{I_c - K_c}{S_c} + 1 & \text{no padding} \\ \left( \frac{I_c - \text{mod}(K_c, 2)}{S_c} \right) + 1 & \text{zero padding} \end{cases} \quad (2.3)$$

adalah jumlah baris dari *output* gambar dimana  $I_r$  and  $I_c$  adalah, berurutan, jumlah baris dan kolom dari input gambar.  $K_r$  and  $K_c$  adalah, berurutan, jumlah baris dan kolom di masing-masing kernel. Gabungan kernel matriksnya berukuran  $C_v \times (O_r$

$\times Oc \times Cx$ ), dan fitur *input* matriks yang menumpuk berukuran  $(Or \times Oc \times Cx) \times (Kr \times Kc)$ [7].

Dengan konversi ini, perhitungan dari titik konvolusi dengan operan  $W, Y$  menjadi[7]:

$$X(Y) \leftarrow Pack(Y) \quad (2.4)$$

$$V(W, Y) \leftarrow WX \quad (2.5)$$

$$\nabla_W^J \leftarrow \nabla_W^J + \nabla_n^J X^T \quad (2.6)$$

$$\nabla_X^J \leftarrow W^T \nabla_n^J \quad (2.7)$$

$$\nabla_Y^J \leftarrow \nabla_Y^J + Unpack(\nabla_X^J) \quad (2.8)$$

Teknik ini memungkinkan pengerjaan paralel untuk operasi matriks yang besar, tetapi menambah pekerjaan di membungkus dan membuka matriks (*pack* dan *unpack*). Dalam kebanyakan kondisi, lebih banyak untungnya daripada kerugiannya. Saat membuat titik konvolusi dengan titik tambahan dan fungsi *element-wise nonlinear*, bias dan *nonlinearity* bisa ditambahkan ke operasi konvolusi[7].

Ada 2 jenis *pooling* yang biasanya digunakan[7]:

- *MaxPooling*: Menerapkan operasi *pooling* maksimum ke *input*  $X$  didalam jendela berukuran  $Kr \times Kc$  untuk setiap kanal. Operasi jendela bergerak mengikuti *subsampling rate* dari *input*  $Sr$  dan  $Sc$

pada arah vertikal dan horizontal secara berurutan. Operasi *pooling* tidak merubah jumlah kanal, sehingga  $C_v = C_x$ .

- *Average Pooling*: Menerapkan rata-rata operasi *pooling* ke *input X* pada jendela berukuran  $K_r \times K_c$  untuk setiap kanal. Operasi jendela bergerak mengikuti *subsampling rate* dari *input S<sub>r</sub>* dan *S<sub>c</sub>* pada arah vertikal dan horizontal secara berurutan.

Arsitektur dari sebuah CNN berisi sebuah pasang lapisan, lapisan konvolusi dan lapisan *pooling*. *Input* berisi angka dari fitur yang dilokalisasi sebagai peta fitur. Ukuran dari peta fitur menjadi kecil pada lapisan atas, seiring dengan semakin banyaknya operasi konvolusi dan *pooling* dilakukan. Biasanya satu atau lebih *hidden layers* yang terhubung ditambahkan di atas lapisan CNN untuk menggabungkan fitur dari semua *frequency bands* sebelum dimasukkan ke *output layer*[1].

## 2.5. Kerangka Berpikir

Penelitian ini dilakukan dengan proses mengumpulkan data terlebih dahulu, lalu melakukan *training* dari data yang sudah dikumpulkan dengan metode *Convolutional Neural Network* yang berada pada *framework* TensorFlow. Lalu membangun aplikasi *speech recognition* dengan menggunakan komputer sebagai perangkat keras untuk menerima input berupa kata yang diucapkan oleh pengguna. Setelah itu melakukan *testing* dan menguji akurasi dari aplikasi yang sudah dibuat.