

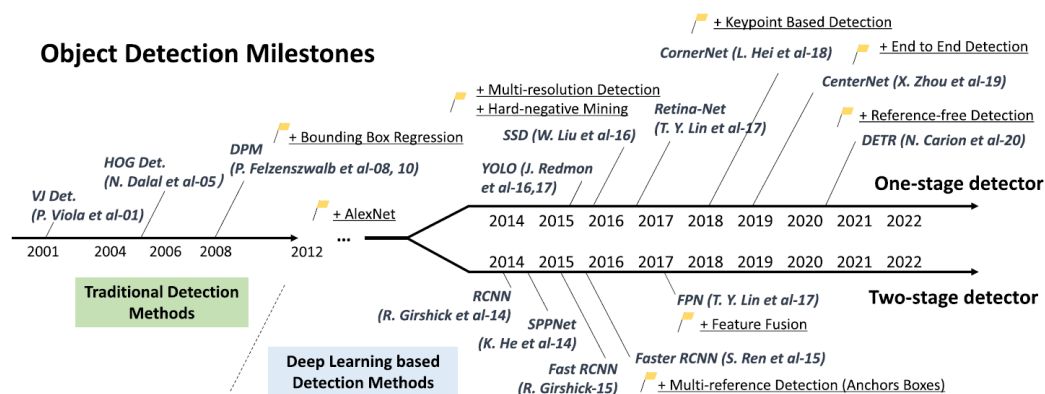
## BAB II

### LANDASAN TEORI

#### 2.1 Deteksi Objek (*Object Detection*)

Deteksi objek merupakan bagian dari visi komputer yang bertugas untuk mendeteksi visual yang memiliki kelas tertentu pada sebuah gambar digital (Zou dkk., 2023). Tujuan utama dari deteksi objek ini adalah mengembangkan sebuah model dan teknik komputasi seperti konvolusi, segmentasi, dan lain-lain untuk mendeteksi objek beserta letak objek tersebut. Dalam deteksi objek terdapat dua komponen utama metrik yang digunakan untuk menghitung performa dari model deteksi objek di antaranya akurasi dan kecepatan deteksi (Zou dkk., 2023).

Deteksi objek pertama kali diawali oleh *Traditional Detector Method* yang dikembangkan oleh (Viola, 2001) dengan nama *Viola Jones Detector (VJ Detection)*. Model yang dikembangkan tersebut menerapkan perhitungan sederhana, yaitu *sliding window* dengan menelusuri semua kemungkinan lokasi dan melakukan penskalaan pada setiap gambar yang ditangkap untuk mengidentifikasi apakah terdapat wajah manusia. Pada masanya, kecepatan deteksi *VJ Detection* paling tinggi jika dibandingkan dengan *detector* lain dengan mengintegrasikan tiga teknik penting, yaitu *integral image*, *feature selection*, dan *detection cascade* (Viola, 2001). Selanjutnya, pada tahun 2012 tercipta suatu inovasi bernama *Deep Convolutional Neural Network (CNN)* (Krizhevsky dkk., 2012) yang menggantikan *Traditional Detector Method* menjadi *Deep Learning Based Detection Method* dari *You Only Look Once (YOLO)* (Zou dkk., 2023).



Gambar 2.1 Perkembangan deteksi objek dalam 20 tahun terakhir (Zou dkk., 2023)

YOLO pertama kali diperkenalkan oleh (Redmon dkk., 2016) dengan menerapkan *single neural network* untuk memproses seluruh gambar secara penuh. Arsitektur jaringan yang digunakan akan membagi gambar menjadi beberapa wilayah dan selanjutnya melakukan prediksi *bounding box* serta menghitung probabilitas untuk setiap kelas pada setiap wilayah secara bersamaan. Model ini dapat meningkatkan kecepatan deteksi dengan pesat, meskipun hasil akurasi deteksi yang dihasilkan harus dikorbankan (Redmon dkk., 2016). Saat ini, versi YOLO telah mencapai versi 8 dengan segala keunggulannya (J. R. Terven & Cordova-esparza, 2024). YOLOv8 dikembangkan oleh (Jocher dkk., 2023) dengan memperbaiki atau meningkatkan arsitektur yang ada pada YOLOv5 (Jocher, 2020). YOLOv5 menggunakan *Cross Stage Partial (CSP) Layer* dibuat untuk meningkatkan kinerja jaringan dengan membagi jalur masukan menjadi dua jalur yang berbeda, kemudian menggabungkan kembali fitur-fitur tersebut (Jocher, 2020). Ini membantu dalam mengatasi masalah degradasi jaringan (*network degradation*) yang bisa terjadi dalam jaringan yang sangat dalam. Pada YOLOv8 *layer* tersebut diganti dengan



## 2.2 *Speech Recognition*

*Speech recognition* merupakan salah satu teknologi yang dapat melakukan konversi suara menjadi sebuah teks yang dapat dikenali oleh mesin (Malik dkk., 2021). Secara ideal *speech recognition* harus dapat menerima input suara yang diberikan oleh manusia, mengenali kata-kata yang diucapkan, dan menggunakan kata-kata tersebut sebagai perintah kepada mesin atau komputer untuk melakukan aksi (Y. dkk., 2017). Perkembangan dari *speech recognition* telah dimulai dari tahun 1950an dan masih berkembang sampai sekarang (Malik dkk., 2021). *Audrey* merupakan sistem *speech recognition* pertama yang dikembangkan oleh laboratorium Bell (Davis dkk., 1952). Sistem yang dirancang hanya dapat mengenali dan membedakan digit-digit berbeda yang diucapkan oleh seorangan pengguna tunggal.

Saat ini, penggunaan *Recurrent Neural Network* (RNN) yang merupakan bagian dari *Long Short Term Memory* (LSTM) menjadi model *speech recognition* paling banyak digunakan dan memiliki kemampuan pengenalan suara yang baik (Malik dkk., 2021). Model tersebut terinspirasi dari penggabungan antara *Hidden Markov Model* (HMM) dengan *Artificial Neural Network* (ANN) (Bourlard & Morgan, 1994). Sehingga model *speech recognition* yang menggunakan RNN dapat mengenali lebih banyak ucapan dalam berbagai aksen karena model tersebut menyimpan pengetahuannya dalam memori jangka pendek pada LSTM. Implementasi penggunaan *speech recognition* yang menggunakan RNN dapat dilakukan pada bahasa pemrograman *python* dengan mudah. Terdapat modul *SpeechRecognition* yang dapat diunduh dan modul tersebut telah mendukung

banyak API untuk *speech recognition* seperti CMU Spinx, *Google Speech Recognition*, *Microsoft Bing Voice Recognition*, dan lain-lain.

### **2.3 Text-to-Speech (TTS)**

*Text-to-Speech* (TTS) atau dapat disebut juga dengan *speech synthesis* merupakan suatu metode yang bertujuan untuk melakukan sintesis ucapan yang dapat dipahami dan dikeluarkan secara alami (Tan dkk., 2021). Secara sederhana, TTS akan melakukan konversi terhadap teks yang menjadi input dengan mengubahnya menjadi sebuah gelombang yang dapat memunculkan suara. Penelitian yang tercatat pertama kali mengenai pengembangan dari TTS ini ada pada tahun 1976 dengan nama model *Articulatory Synthesis* (Coker, 1976). *Articulatory Synthesis* adalah pendekatan dalam pembuatan sistem sintesis ucapan yang mencoba meniru bagaimana manusia menghasilkan suara dengan memodelkan gerakan fisik dari organ-organ bicara seperti lidah dan bibir (Coker, 1976). Dengan menggunakan model ini, sistem dapat menghasilkan suara yang lebih realistis yang dapat disesuaikan dengan berbagai gaya bicara dan aksen pada masanya.

Dengan kemajuan dalam *deep learning*, diperkenalkan TTS berbasis jaringan saraf (*neural TTS*), yang menggunakan *neural network* sebagai kerangka model untuk sintesis suara. WaveNet (Oord dkk., 2016) kemudian diusulkan untuk langsung menghasilkan gelombang suara dari fitur linguistik yang dianggap sebagai model TTS berbasis jaringan saraf modern pertama (Tan dkk., 2021). Model lain seperti DeepVoice (Arik dkk., 2017) masih mengikuti komponen-komponen dalam sintesis parametrik statistik, tetapi diperbarui dengan model berbasis jaringan saraf

yang sesuai. Selain itu, beberapa model *end-to-end* seperti Tacotron (Wang dkk., 2017) , Deep Voice 3 (Ping dkk., 2018), dan FastSpeech (Ren dkk., 2019) menyederhanakan modul analisis teks dan langsung menggunakan urutan karakter/fonem sebagai input, serta menyederhanakan fitur akustik dengan *mel-spectrogram*. Kemudian, sistem TTS *end-to-end* sepenuhnya dikembangkan untuk menghasilkan gelombang suara langsung dari teks, seperti ClariNet (Ping dkk., 2019), FastSpeech 2 (Ren dkk., 2021), dan *End-to-End Adversarial Text-to-Speech* (EATS) (Donahue dkk., 2021). Dibandingkan dengan sistem TTS sebelumnya, kelebihan *neural text-to-speech* meliputi kualitas suara yang tinggi dalam hal kejelasan dan kealamiannya, serta persyaratan yang lebih sedikit dalam pra-pemrosesan manusia dan pengembangan fitur.

#### **2.4 Penelitian Terkait**

Saat ini telah banyak pengembangan-pengembangan terhadap perangkat yang digunakan untuk membantu penderita tunanetra dalam menjalani aktivitasnya sehari-sehari. Pengembangan yang dilakukan tidak semata-mata hanya dilakukan tanpa aspek-aspek yang mempengaruhi kualitas perangkat. Terdapat beberapa aspek yang menjadi fitur utama dalam mengukur performa pengembangan perangkat yang diantaranya *capturing device*, *working hour*, *response time*, *coverage area*, *feedback*, *working range*, *weight*, *robustness*, dan *cost* (Mashiata dkk., 2022). Pengukuran performa tersebut didasarkan pada tantangan-tantangan yang dialami oleh para penderita tunanetra serta berasal dari lingkungan para penderita, kehidupan sosialnya, kesulitan dalam menggunakan teknologi dan lain sebagainya (Manjari dkk., 2020).

Pemanfaatan *object detection* dan *text-to-speech feedback* berbasis *Artificial Intelligence* (AI) semakin banyak digunakan untuk mengatasi setiap tantangan yang ada dan memenuhi setiap aspek yang dibutuhkan (Walle dkk., 2022). Namun, terdapat masalah pada keakuratan model deteksi dan *feedback* yang diberikan oleh model apakah memiliki kesesuaian dengan apa yang dibutuhkan oleh pengguna. Sehingga, untuk menjawab permasalahan tersebut dibuatlah *state of the art* terkait dengan penelitian yang dilakukan serta melihat kesesuaian pemanfaatan dari *object detection* dan TTS pada perangkat pembantu tunanetra.

Tabel 2.1 *State of the art* penelitian terkait

Penulis	Judul	Metode Deteksi Objek	Metode TTS	Keterangan
(Guravaiah dkk., 2022)	<i>Third Eye: Object Recognition and Speech Generation for Visually Impaired</i>	YOLOv5	<i>google Text-to-Speech</i> (gTTS) dan <i>pyttsx3</i>	<b>Kelebihan</b> : Hasil perhitungan matrik mAP menunjukkan nilai mAP:0.5 pada 39% dan mAP:0.95 pada 25%, <b>Kekurangan</b> : Hasil mAP yang dihasilkan oleh model YOLOv5 masih kurang baik karena untuk hasil mAP:0.5 saja berada di bawah 50%.
(Setiadi dkk., 2020)	<i>Navigation and Object Detection for Blind Persons Based on Neural Network</i>	CNN dan LiDAR ( <i>Light Detection and Ranging</i> )	<i>Voice attitudes information</i>	<b>Kelebihan</b> : menunjukkan akurasi hasil pendeteksian jalur pejalan kaki sebesar 89.7% pada lux<15000 dan 87.5% untuk akurasi <i>object detection</i> <b>Kekurangan</b> : Penelitian yang dilakukan tidak menjelaskan bagaimana arsitektur yang digunakan dan tidak memperlihatkan bagaimana optimasi model bekerja
(Ravindra Karmarkar &	<i>Object Detection System for The</i>	YOLO	gTTS	<b>Hasil</b> : Menunjukkan bahwa model dapat

Honmane, 2021)	<i>Blind with Voice Guidance</i>			melakukan pendeteksian objek dan mengeluarkan suara dari hasil pendeteksian. <b>Kekurangan</b> : Penelitian tersebut tidak menunjukkan perbandingan dan akurasi hasil model
(Najm dkk., 2022)	<i>Assisting Blind People Using Object Detection with Vocal Feedback</i>	YOLOv3	gTTS	<b>Hasil</b> : Untuk gambar 608x608 memiliki akurasi sebesar 90.17%, untuk gambar 512x512 memiliki akurasi sebesar 87.32%, dan untuk gambar 416x416 memiliki akurasi sebesar 97.17%. <b>Kekurangan</b> : Model deteksi objek yang digunakan hanya terpakai pada YOLOv3 saja
(Wong dkk., 2019)	<i>Convolutional Neural Network for Object Detection System for Blind People</i>	CNN dengan arsitektur SSD Mobilenet	<i>Voice synthesis</i>	<b>Hasil</b> : Model menunjukkan akurasi sebesar 73.7% pada <i>dataset</i> CIFAR-10. <b>Kekurangan</b> : Model hanya dapat mencapai FPS sebesar 5 yang membuat pendeteksian secara real time menjadi kurang baik.
(Abdurrasyid dkk., 2019)	<i>Detection of immovable objects on visually impaired people walking aids</i>	<i>Template matching</i>	-	<b>Hasil</b> : Hasil menunjukkan bahwa ketika kamera diletakkan pada 90° dengan objek, hasil deteksi menjadi lebih maksimal karena noise yang didapatkan sedikit. <b>Kekurangan</b> : Ketika kamera berada pada 45° dan 135° dengan objek, hasil deteksi menjadi tidak maksimal dan mendapatkan akurasi terendah pada pengujian.
(Konaite dkk., 2021)	<i>Smart Hat for the blind with Real-Time Object Detection using Raspberry Pi and TensorFlow Lite</i>	SSD MobileNet v2 320x320	SPEAK	<b>Hasil</b> : Deteksi <i>real-time</i> FPS sebesar 5.26 frame/s pada raspberry pi 4. Untuk waktu deteksi gambar



				mencapai waktu tercepat, yaitu 19 ms. <b>Kekurangan</b> : Model memiliki nilai mAP yang paling kecil jika dibandingkan dengan model lain yang dibandingkan yaitu 20.2%
(Mahendru & Dubey, 2021)	<i>Real Time Object Detection with Audio Feedback using Yolo vs. Yolo_v3</i>	YOLOv3 dengan Darknet53 dan YOLO dengan SSD MobilNet	gTTS	<b>Kelebihan</b> : Pendeteksian <i>multi-class</i> pada satu gambar dengan baik. rata-rata akurasi hasil deteksi gambar yang dilakukan oleh model ada pada 65-98 % <b>Kekurangan</b> : Penelitian ini masih menggunakan YOLOv3 yang original sehingga tidak membangkitkan performa yang maksimal bagi model
(Ramalakshmi dkk., 2020)	<i>Object Detector for Visually Impaired with Distance Calculation for Humans</i>	Integrasi YOLO dengan <i>Haar Classifier</i>	gTTS	<b>Kelebihan</b> : Model dapat melakukan deteksi pada objek yang saling berhimpitan. Selain itu, model dapat melakukan kalkulasi perkiraan jarak antara objek manusia dengan kamera melalui <i>Haar classifier</i> <b>Kekurangan</b> : Penelitian ini tidak memperlihatkan bagaimana hasil performa model secara menyeluruh
(Alzahrani & Al-Baity, 2023)	<i>Object Recognition System for the Visually Impaired: A Deep Learning Approach using Arabic Annotation</i>	Mask R-CNN	gTTS	<b>Kelebihan</b> : Model Mask R-CNN memiliki nilai mAP sebesar 83.9% dengan objek gambar yang berhasil dideteksi beserta letak objek tersebut. <b>Kekurangan</b> : Terdapat beberapa kelas objek yang memiliki nilai mAP rendah dan model ini masih belum mendukung pendeteksian secara <i>real time</i> .

(R. C. Joshi dkk., 2020)	<i>Efficient Multi-Object Detection and Smart Navigation Using Artificial Intelligence for Visually Impaired People</i>	YOLOv3 dengan tambahan sensor ultrasonic	<i>Optical Character Recognizer (OCR)</i>	<p><b>Kelebihan</b> : 95.19% untuk akurasi deteksi dan 99.69% untuk akurasi pengenalan objek.</p> <p><b>Kekurangan</b> : Dibutuhkan tambahan kelas pada dataset untuk memaksimalkan pendeteksian pada model</p>
(Alahmadi dkk., 2023)	<i>Enhancing Object Detection for VIPs Using YOLOv4_Resnet101 and Text-to-Speech Conversion Model</i>	YOLOv4 dengan backbone Resnet101	pyttsx3	<p><b>Kelebihan</b> : Akurasi model sebesar 96.34% dan error rate 0.073%. Kecepatan FPS yang dihasilkan ada pada 60-80 frame/s untuk input gambar 416x416</p> <p><b>Kekurangan</b> : Memerlukan <i>resource</i> yang besar dalam proses pelatihan maupun percobaan.</p>
(Noman dkk., 2020)	<i>Portable offline indoor object recognition system for the visually impaired</i>	Faster R-CNN	eSpeakNG	<p><b>Kelebihan</b> : Akurasi model sebesar 73.24% dengan kesalahan pendeteksian jarak sekitar 5%.</p> <p><b>Kekurangan</b> : Akurasi yang dihasilkan masih kurang karena model yang digunakan hanya menggunakan Tensorflow Object Detection API saja</p>
(Lee & Cho, 2022)	<i>Automatic Object Detection Algorithm-Based Braille Image Generation System for the Recognition of Real-Life Obstacles for Visually Impaired People</i>	YOLOv3	-	<p><b>Kelebihan</b> : Hasilnya, sistem bantuan hidup yang diusulkan dalam penelitian ini terbukti efisien dan berguna dengan akurasi rata-rata konversi braille 85%, akurasi deteksi 90% atau lebih, dan waktu konversi braille rata-rata 6,6 detik</p> <p><b>Kekurangan</b> : Hasil ekstraksi objek yang diperoleh melalui GrabCut menggunakan koordinat objek yang terdeteksi dengan YOLOv3 tidak sesuai dengan objek sebenarnya.</p>

(Rahman & Sadi, 2021)	<i>IoT Enabled Automated Object Recognition for the Visually Impaired</i>	SSD MobilNet	-	<p><b>Kelebihan</b> : Akurasi keseluruhan sistem yang diusulkan dalam deteksi dan pengenalan objek masing-masing adalah 99,31% dan 98,43%.</p> <p><b>Kekurangan</b> : Penelitian tidak menjelaskan metode TTS apa yang digunakan.</p>
(Rocha dkk., 2023)	<i>Using Object Detection Technology to Identify Defects in Clothing for Blind People</i>	YOLOv5	<i>Smartphone Audio Feedback</i>	<p><b>Kelebihan</b> : Dari perbandingan yang dilakukan model YOLOv5s6 dapat melakukan waktu inferensi yang paling cepat, yaitu 0.0092s, sedangkan untuk nilai mAP terbaik ada pada model YOLOv5l6 dengan mAP0.5 sebesar 76%.</p> <p><b>Kekurangan</b> : Metode optimasinya masih menggunakan SGD yang masih kurang maksimal jika dibandingkan dengan metode lain</p>
(Kadhim & Oleiwi, 2022)	<i>Blind Assistive System Based on Real Time Object Recognition using Machine Learning</i>	YOLOv3	-	<p><b>Kelebihan</b> : Hasil menunjukkan bahwa YOLOv3 dapat melakukan deteksi dengan baik pada kelas-kelas yang ada pada dataset.</p> <p><b>Kekurangan</b> : Tidak adanya perbandingan dengan penelitian lain membuat model pada penelitian ini tidak dapat dipastikan performanya</p>
(Yohannes dkk., 2020)	<i>Robot Eye: Automatic Object Detection And Recognition Using Deep Attention Network to Assist Blind People</i>	<i>Deep Attention Network</i>	-	<p><b>Kelebihan</b> : Metode yang diusulkan mencapai tingkat akurasi sekitar 81%, lebih baik daripada YOLO v3 original.</p> <p><b>Kekurangan</b> : Performa model yang dihasilkan tidak mencapai performa yang maksimal karena tidak adanya proses optimasi model</p>

(Islam dkk., 2023)	<i>Deep learning based object detection and surrounding environment description for visually impaired people</i>	TensorFlow Object Detection API dan SSDLite MobileNetV2	Google text-to-speech, PyAudio, playsound	<p><b>Kelebihan</b> : Akurasi yang dihasilkan oleh model ada pada 88.89% dan untuk FPS-nya berada pada 2.15 frame/s</p> <p><b>Kekurangan</b> : Ukuran kumpulan data cuaca yang digunakan untuk mengembangkan ambience mode relatif kecil.</p>
(Khan dkk., 2023)	<i>Outdoor mobility aid for people with visual impairment: Obstacle detection and responsive framework for the scene perception during the outdoor mobility of people with visual impairment</i>	YOLOv5 dan Mask R-CNN	Google services	<p><b>Kelebihan</b> : Hasil pengukuran mAP(0.5) menunjukkan untuk YOLOv5s mendapatkan nilai 97% dan untuk Mask R-CNN mendapat nilai 93%.</p> <p><b>Kekurangan</b> : Kelas yang digunakan pada dataset terbilang sedikit karena hanya terdiri dari 8 kelas saja.</p>
(Mukhiddinov & Cho, 2021)	<i>Smart Glass System Using Deep Learning for the Blind and Visually Impaired</i>	DETR (Detection Transformer)	OCR	<p><b>Kelebihan</b> : Model memiliki nilai mAP:0.5 sebesar 63.5%. Untuk model ekstraksi objek menonjol menghasilkan nilai maxFM, MAE, dan WFM secara berturut-turut adalah 0.814, 0.058, dan 0.725.</p> <p><b>Kekurangan</b> : Dalam beberapa situasi, model deteksi objek mendeteksi lebih dari sepuluh objek, di mana beberapa di antaranya adalah objek kecil atau terdeteksi dengan tidak benar,</p>
(Mandhala dkk., 2020)	<i>Object Detection Using Machine Learning for Visually Impaired People</i>	Retina Net, YOLOv3, dan Tiny-YOLOv3	-	<p><b>Kelebihan</b> : Akurasi masing-masing untuk Retina Net 80.5%, YOLOv3 71.1%, dan Tiny-YOLOv3 76.3%.</p> <p><b>Kekurangan</b> : Tidak adanya proses optimasi membuat model tidak menghasilkan performa yang maksimal.</p>
(Pi dkk., 2022)	<i>Low-Cost Smart Glasses for Blind Individuals using Raspberry Pi 2</i>	YOLOv3	OCR	<p><b>Kelebihan</b> : Dengan jarak kurang lebih 5 meter, model dengan</p>

				<p>baik dapat mengenali objek yang ditangkap.</p> <p><b>Kekurangan</b> : Terdapat perbedaan dalam pendeteksian beberapa objek karena jarak yang menyebabkan pendeteksiannya tidak dapat mendeteksi hal yang spesifik pada percobaan tertentu</p>
(Hsieh dkk., 2021)	<i>A CNN-Based Wearable Assistive System for Visually Impaired People Walking Outdoors</i>	FAST-SCNN dan YOLOv5s	<i>Voice prompt manager</i>	<p><b>Kelebihan</b> : Performa dari model hasil pelatihan Fast-SCNN akurasi pixel 95%. Sedangkan untuk model deteksi rintangan yang menggunakan YOLOv5s menunjukkan performa mAP50:95 sebesar 52.6%, mAP50 sebesar 85.1%, dan mAP75 sebesar 58.5%</p> <p><b>Kekurangan</b> : Kecepatan deteksi yang masih kurang dapat membahayakan penggunaa jika tiba-tiba muncul objek di depannya</p>
(Ganesan dkk., 2022)	<i>Deep Learning Reader for Visually Impaired</i>	CNN dan LSTM	Text-to-Speech API	<p><b>Kelebihan</b> : Dengan menggunakan Resnet arsitektur performa model menunjukkan hasil akurasi sebesar 83%. Selain itu, sistem dapat melakukan pemahaman terhadap gambar yang dideteksi dengan menambahkan caption pada setiap hasil deteksi yang membuatnya semakin natural.</p> <p><b>Kekurangan</b> : Tidak ada fitur untuk menghitung jarak objek</p>
(Chou dkk., 2023)	<i>A Lightweight Robust Distance Estimation Method for Navigation Aiding in Unsupervised Environment Using Monocular Camera</i>	YOLOv4	-	<p><b>Kelebihan</b> : Dengan baik model dapat menghitung jarak objek pada rentang di bawah 10 meter dengan total kesalahan sebesar 4%. Selain itu, penggunaan YOLOv4 membuat</p>

				<p>proses deteksi objek menjadi lebih baik lagi.  <b>Kekurangan</b> : Saat ini, model hanya dapat mengenali dua jenis objek publik di luar ruangan, yang berarti model hanya dapat berfungsi di luar ruangan. Selain itu, benda bergerak seperti manusia dan mobil juga belum dapat terdeteksi</p>
--	--	--	--	--

Penggunaan YOLO untuk model deteksi objek sebagai komponen dari alat pembantu tunanetra menunjukkan kesuksesan yang sangat luar biasa. Dibuktikan dengan beberapa penelitian yang menggunakan model tersebut sebagai basis deteksi objeknya pada Tabel 2.1. Selain itu, penggunaan gTTS masif digunakan sebagai bentuk implementasi untuk TTS pada keluaran model. Dari *state of the art* yang ada dan pernyataan yang diberikan, terciptalah ide penelitian pada proposal ini dengan mengusung model yang memiliki kemampuan deteksi objek, translasi TTS, *speech recognition*, dan menghitung estimasi jarak antara kamera dengan objek. Tabel 2.2 menunjukkan perbandingan capaian yang didapatkan pada penelitian ini dengan capaian yang telah didapatkan oleh penelitian sebelumnya. Perbandingan tersebut mengacu pada dukungan model versi YOLO, *parameter optimization*, *voice feedback*, *speech recognition*, dan *distance estimation*

Tabel 2.2 Perbandingan target capaian penelitian dengan penelitian terkait

Penelitian	YOLO	<i>Parameter Optimization</i>	<i>Voice Feedback</i>	<i>Speech Recognition</i>	<i>Distance Estimation</i>
(Guravaiah dkk., 2022)	YOLOv5	-	√	-	-
(Ravindra Karmarkar &	YOLOv1	-	√	-	-

Honmane, 2021)					
(Najm dkk., 2022)	YOLOv3	-	√	-	-
(Mahendru & Dubey, 2021)	YOLOv3	-	√	-	-
(Ramalakshmi dkk., 2020)	YOLOv1	-	√	-	√
(N. Joshi dkk., 2021)	YOLOv3	-	√	-	-
(Alahmadi dkk., 2023)	YOLOv4	-	√	-	-
(Lee & Cho, 2022)	YOLOv3	-	√	-	-
(Kadhim & Oleiwi, 2022)	YOLOv3	-	√	-	-
(Nugraha, 2024)	YOLOv8	√	√	√	√

Tabel 2.2 menunjukkan perbandingan capaian setiap penelitian terkait dengan penelitian yang dilakukan. Pemilihan penelitian terkait pada tabel tersebut didasari pada penggunaan model YOLO yang digunakan sebagai model objek deteksi untuk membantu penyandang tunanetra. Pada penelitian (Guravaiah dkk., 2022) versi YOLO yang digunakan adalah versi 5 dengan dikombinasikan oleh TTS menggunakan gTTS dan pyttsx3 sehingga menghasilkan suatu sistem bernama Third Eye. Selain itu, penambahan 15 kelas baru pada MS COCO *dataset* membuat model yang dihasilkan dapat lebih banyak mengenali banyak objek. (Ramalakshmi dkk., 2020; Ravindra Karmarkar & Honmane, 2021) menggunakan YOLOv1 sebagai dasar dari model objek deteksinya. Meskipun menggunakan YOLO versi pertama, tetapi model tersebut dapat bekerja dengan baik dalam melakukan deteksi pada kelas-kelas yang terdaftar pada *dataset* COCO. Pada penelitian (N. Joshi dkk., 2021; Kadhim & Oleiwi, 2022; Lee & Cho, 2022; Mahendru & Dubey, 2021; Najm

dkk., 2022) model deteksi yang digunakan adalah YOLOv3 dimana versi ini terkenal dengan fleksibilitas dan ketangguhannya dalam hal deteksi objek. Hal ini membuktikan bahwa penggunaan YOLOv3 saat ini masih relevan untuk digunakan. Penambahan fitur feedback berupa text-to-speech menjadi pelengkap dari semua sistem yang dibangun pada penelitian-penelitian tersebut. (Alahmadi dkk., 2023) mengembangkan YOLOv4 dengan mengganti *backbone*-nya dari Darknet menjadi Resnet101. Sehingga dihasilkan sebuah model deteksi yang memiliki performa yang sangat baik dengan rata-rata akurasi sebesar 96.34% pada 16 kelas yang ada pada MS COCO *dataset*.

Penelitian-penelitian terkait tersebut memiliki kekurangan yang identik, yaitu tidak dilakukannya pemilihan parameter yang baik khususnya parameter untuk masalah optimasi. Selain itu, hanya segelintir penelitian terkait yang menggunakan fitur *distance estimation* bahkan tidak ada sama sekali yang menggunakan fitur *speech recognition*. Padahal kedua fitur tersebut dapat meningkatkan kebergunaan dari sistem yang dihasilkan bagi para penyandang tunanetra. Maka dari itu, penelitian ini menghasilkan kebaruan yang tidak dimiliki oleh penelitian-penelitian terkait pada Tabel 2.2. Penelitian ini menerapkan metode *parameter optimization* menggunakan PSO pada pemilihan nilai *learning rate* dan *beta 1 Adam optimization* yang paling optimal bagi proses pelatihan model deteksi. Selain itu, penambahan fitur TTS, *speech recognition*, dan *distance estimation* menggunakan perhitungan *monocular camera* ditambahkan untuk mendukung kegunaan sistem bagi para penyandang tunanetra.