

---

## IMAGE TO TEXT TO SPEECH CONVERTER

---

**\*Lalit Kumar, Navneet Yadav, Vinit, Arun Kumar**

---

Department of Computer Science & Engineering, R D Engineering College Ghaziabad, UP,  
India-201206.

**Article Received: 27 March 2026, Article Revised: 17 April 2026, Published on: 07 May 2026**

**\*Corresponding Author: Lalit Kumar**

Department of Computer Science & Engineering, R D Engineering College Ghaziabad, UP, India-201206.

DOI: <https://doi-doi.org/101555/ijarp.4275>

### ABSTRACT

This research paper presents the design, development, and evaluation of an AI-powered web application that integrates Optical Character Recognition (OCR) technology with browser-based Text-to-Speech (TTS) synthesis to create an accessible, full-stack platform for extracting spoken content from digital images. Named the Image-to-Text-to-Speech (ITIS) Converter, this system addresses a critical accessibility gap encountered by visually impaired individuals, language learners, and professionals dealing with large volumes of scanned or photographed textual content.

The system employs a React-based frontend for an intuitive user experience, Node.js and Express for backend API management, and Mongo DB as a persistent document store for saving and retrieving extracted stories. The core OCR functionality is powered by Tesseract.js — the JavaScript port of Google's Tesseract OCR engine — which operates entirely in the browser through Web Assembly, eliminating the need for server-side image processing. Speech synthesis is handled by the native Web Speech API, offering real-time customization of voice, rate, and pitch without additional server load.

Through systematic evaluation across five image categories — printed text, handwritten content, low-resolution scans, multilingual documents, and natural scene text — the application achieved an average character accuracy of 91.4% under optimal conditions. Key findings indicate that image preprocessing and lighting quality are the primary determinants of OCR accuracy. The system's modular architecture and separation of frontend and backend concerns make it extensible toward future enhancements including authentication, audio download, and multi-language support.

This paper documents the full system lifecycle including literature review, system architecture, implementation details, testing methodologies, performance evaluation, limitations, and a proposed roadmap for future development. The project demonstrates that modern browser APIs and JavaScript libraries are sufficiently mature to power production-grade accessibility tools without dependency on heavyweight server-side AI infrastructure.

**KEYWORDS:** Optical Character Recognition (OCR), Text-to-Speech (TTS), Tesseract.js, Web Speech API, React, Node.js, MongoDB, Accessibility, Full-Stack Web Application, Web Assembly

## 1. INTRODUCTION

In the modern digital landscape, information is increasingly disseminated through image-based formats — scanned documents, photographed text, screenshots, and digital posters. While these formats are visually rich, they present significant accessibility barriers for individuals with visual impairments, reading disabilities, or language acquisition challenges. According to the World Health Organization (WHO), over 2.2 billion people globally experience some form of vision impairment, and a substantial proportion of this population relies on assistive technology to consume textual information.

The convergence of these technologies presents an opportunity to create a seamless, low-cost accessibility tool that enables users to photograph or upload any image containing text and immediately receive both a digital transcript and an audio rendering of that text. This project, the Image-to-Text-to-Speech (ITIS) Converter, was conceived and built in response to this opportunity, providing a practical demonstration of how modern web technologies can be orchestrated to serve real accessibility needs.

### 1.1 Problem Statement

Despite the proliferation of digital devices capable of capturing high-resolution images, there remains a significant gap between the generation of image-based textual content and its accessibility to individuals who cannot visually process such images. Existing solutions often require paid subscriptions, desktop applications, or specialized hardware, creating economic and logistical barriers. Furthermore, most web-based OCR tools do not integrate speech synthesis, requiring users to perform multiple discrete steps — OCR, copying text, switching to a TTS tool — to achieve the same outcome that an integrated platform could deliver in a single workflow.

This project addresses the following core problem: How can a unified, browser-based web application leverage modern JavaScript libraries and APIs to extract text from images and deliver that text as synthesized speech, while also providing persistent storage and retrieval capabilities, in a manner that is accessible, efficient, and extensible?

## **1.2 Optical Character Recognition: Historical Overview**

Optical Character Recognition is the electronic or mechanical conversion of images of typed, handwritten, or printed text into machine-encoded text. The concept was first patented by Emanuel Goldberg in 1931, who developed a machine that could read characters and convert them to telegraph code. Commercial OCR systems began to emerge in the 1950s and 1960s, primarily used by banks and postal services for reading standardized printed characters.

The field advanced significantly with the introduction of neural network-based approaches in the 1990s. Yann LeCun's landmark work on Convolutional Neural Networks (CNNs) and gradient-based learning for document recognition (1998) laid the groundwork for modern deep learning-based OCR systems. Google's Tesseract, originally developed at Hewlett-Packard in the 1980s and subsequently open-sourced by Google in 2005, became one of the most widely used OCR engines due to its accuracy, multilingual support, and extensibility.

Smith (2007) documented the architecture of Tesseract 2.0, highlighting its use of an adaptive classifier that improves accuracy on a per-document basis. The transition from traditional image processing pipelines to LSTM (Long Short-Term Memory) neural networks in Tesseract 4.0 (2018) represented a paradigm shift, improving accuracy by 30-40% on complex documents compared to the previous version. The JavaScript port, Tesseract.js, created by Jerome Wu and maintained by the open-source community, brings this capability to the browser environment via WebAssembly compilation.

Breuel et al. (2013) demonstrated that LSTM-based sequence recognition models significantly outperform traditional HMM-based approaches for degraded historical document recognition. Their work on OCRopus, an open-source OCR pipeline, influenced the development of Tesseract's neural network mode. Contemporary research by Shi et al. (2016) introduced the End-to-End Trainable Neural Network for Scene Text Recognition (CRNN), which combines CNN feature extraction with LSTM sequence modeling — an architecture that underlies many modern OCR tools.

## **1.3. Text-to-Speech Technology**

Text-to-Speech synthesis has undergone a dramatic evolution from early formant synthesis systems of the 1960s to today's neural network-based voices that are nearly indistinguishable from human speech. The first intelligible TTS systems, such as the Bell Labs Vocoder and

Dennis Klatt's DECTalk, used rule-based formant synthesis that produced robotic-sounding output.

- **Backend Technologies**
- **Node.js**
- **Express.js**
- **CORS Middleware**
- **Dotenv**

## **2. Database Technologies**

### **2.1 Mongo DB**

Mongo DB is a cross-platform, document-oriented NoSQL database that stores data in flexible, JSON-like BSON (Binary JSON) documents. Unlike relational databases that require a predefined schema, MongoDB allows documents within a collection to have different fields, making it highly adaptable to evolving data requirements. For the ITIS Converter, MongoDB's document model naturally accommodates story objects with variable-length text content.

### **2.2 Mongoose**

Mongoose is an Object Data Modeling (ODM) library for MongoDB and Node.js. It provides a schema-based solution for modeling application data, including built-in type casting, validation, query building, and middleware (pre/post hooks). The ITIS Converter uses Mongoose to define the Story schema, enforce required field validation, and provide a clean API for database operations.

## **3. METHODOLOGY**

### **3.1 Project Setup and Configuration**

The project follows a monorepo-style directory organization with separate backend and frontend subdirectories. This structure allows the two components to be developed and deployed independently while residing in the same version control repository.

#### **3.1.1 Backend Initialization**

The backend is initialized using npm init with the following key dependencies defined in package.json:

```
{ "dependencies": {  
  "express": "^4.18.2",  
  "mongoose": "^7.5.0",
```

```
"cors":    "^2.8.5",
"dotenv":  "^16.3.1"
},
"scripts": { "start": "node server.js", "dev": "nodemon server.js" }
}
```

### 3.1.2 Frontend Initialization

The frontend is bootstrapped with Create React App and configured with the proxy setting to route API calls to the backend server during development:

```
// package.json (frontend)
"proxy": "http://localhost:5000"
```

This proxy configuration allows the frontend to make API calls using relative paths (e.g., /api/save) without CORS issues during development, as the CRA development server forwards these requests to the Express server.

## 3.2 Backend Implementation

### 3.2.1 Server Initialization and Middleware

The Express server is initialized with the following middleware stack, applied in order:

```
const express = require('express');
const mongoose = require('mongoose');
const cors = require('cors');
require('dotenv').config();

const app = express();
app.use(cors());           // Cross-origin support
app.use(express.json());   // JSON body parsing
app.use(express.urlencoded({ extended: true }));
```

### 3.2.2 MongoDB Connection

The MongoDB connection is established asynchronously using Mongoose's connect() method, with error handling and a connection success callback:

```
mongoose.connect(
  process.env.MONGO_URI || 'mongodb://localhost:27017/speechApp',
```

```
{ useNewUrlParser: true, useUnifiedTopology: true }
)
.then(() => console.log('MongoDB connected successfully'))
.catch(err => {
  console.error('MongoDB connection error:', err);
  process.exit(1);
});
```

The connection string defaults to a local MongoDB instance, allowing the application to run without external dependencies in development. The `MONGO_URI` environment variable can be set to a MongoDB Atlas connection string for production deployment.

### 3.2.3 Story Model

The Mongoose schema and model are defined as follows, incorporating validation and computed metadata fields:

```
const StorySchema = new mongoose.Schema({
  title: {
    type: String,
    required: [true, 'Story title is required'],
    trim: true,
    maxlength: [200, 'Title cannot exceed 200 characters']
  },
  content: {
    type: String,
    required: [true, 'Story content is required']
  },
  date: { type: Date, default: Date.now },
  wordCount: Number,
  charCount: Number
});
const Story = mongoose.model('Story', StorySchema);
```

### 3.2.4 Save Endpoint

The POST `/api/save` endpoint validates the incoming payload, computes word and character counts, persists the document, and returns the created document's ID:

```
app.post('/api/save', async (req, res) => {
  try {
    const { title, content } = req.body;
    if (!title || !content) {
      return res.status(400).json({
        success: false, error: 'Title and content are required'
      });
    }
    const story = new Story({
      title, content,
      wordCount: content.split(/\s+/).length,
      charCount: content.length
    });
    await story.save();
    res.status(201).json({ success: true, id: story._id });
  } catch (err) {
    res.status(500).json({ success: false, error: err.message });
  }
});
```

### 3.2.5 Get Stories Endpoint

The GET `/api/stories` endpoint retrieves all stored stories, sorted by creation date in descending order (most recent first):

```
app.get('/api/stories', async (req, res) => {
  try {
    const stories = await Story
      .find({}, 'title date wordCount charCount')
      .sort({ date: -1 })
      .limit(50);
    res.json({ success: true, stories });
  } catch (err) {
```

```

    res.status(500).json({ success: false, error: err.message });
  }
});

```

### 3.3 Frontend Implementation

#### 3.3.1 Application State

The App.js component manages the following state variables using React's useState hook:

| State Variable | Type                 | Initial Value | Purpose                      |
|----------------|----------------------|---------------|------------------------------|
| image          | File   null          | null          | Selected image file          |
| imageUrl       | string               | "             | Object URL for image preview |
| extractedText  | string               | "             | OCR result text              |
| isProcessing   | boolean              | FALSE         | OCR in-progress flag         |
| progress       | number               | 0             | OCR progress (0–100)         |
| isSpeaking     | boolean              | FALSE         | TTS active flag              |
| selectedVoice  | SpeechSynthesisVoice | null          | Active TTS voice             |
| voices         | Array                | []            | Available browser voices     |
| rate           | number               | 1             | TTS speech rate              |
| pitch          | number               | 1             | TTS pitch                    |
| storyTitle     | string               | "             | Title for saving             |
| savedStories   | Array                | []            | Fetches stories list         |

#### 3.3.2 Component Lifecycle

The use Effect hook is used for two side effects in the main component:

1. Fetching the list of saved stories from the backend on initial mount.
2. Listening for the voiceschanged event on the speechSynthesis object, which fires asynchronously when the browser's voice list is populated.

```

useEffect(() => {
  fetchStories();
  const loadVoices = () => {
    const v = window.speechSynthesis.getVoices();
    setVoices(v);
    const google = v.find(x => x.name.includes('Google'));
    const hindi = v.find(x => x.lang === 'hi-IN');
    setSelectedVoice(google || hindi || v[0] || null);
  };
  window.speechSynthesis.onvoiceschanged = loadVoices;
  loadVoices();
}, []);

```

### 3.4 OCR Module Implementation

#### 3.4.1 Tesseract.js Integration

The OCR pipeline is implemented as an asynchronous function that is called when a new image is selected. The function creates a Tesseract worker, initializes it with English language data, and calls the recognize() method with a progress callback that updates the UI:

```
import Tesseract from 'tesseract.js';

const handleOCR = async (file) => {
  setIsProcessing(true);
  setProgress(0);
  setExtractedText("");
  try {
    const worker = await Tesseract.createWorker({
      logger: m => {
        if (m.status === 'recognizing text') {
          setProgress(Math.floor(m.progress * 100));
        }
      }
    });
    await worker.loadLanguage('eng');
    await worker.initialize('eng');
    const { data } = await worker.recognize(file);
    setExtractedText(data.text.trim());
    await worker.terminate();
  } catch (err) {
    console.error('OCR Error:', err);
    setExtractedText('Error during OCR. Please try another image.');
```

### 3.4.2 Image Preprocessing Considerations

Tesseract.js performs best on high-contrast, cleanly segmented text images. The current implementation passes raw image data directly to the OCR engine without preprocessing. Future versions may incorporate Canvas API-based preprocessing to convert images to grayscale, increase contrast using histogram equalization, and apply thresholding to binarize the image before OCR processing, which can improve recognition accuracy on low-quality scans by 15–25%.

## 3.5 Text-to-Speech Module Implementation

### 3.5.1 Speech Synthesis

The TTS playback is initiated by the `handleSpeak` function, which creates and configures a `SpeechSynthesisUtterance` object:

```
const handleSpeak = () => {
  if (!extractedText) return;
  window.speechSynthesis.cancel(); // Stop any ongoing speech
  const utterance = new SpeechSynthesisUtterance(extractedText);
  if (selectedVoice) utterance.voice = selectedVoice;
  utterance.rate = rate;
  utterance.pitch = pitch;
  utterance.onstart = () => setIsSpeaking(true);
  utterance.onend = () => setIsSpeaking(false);
  utterance.onerror = (e) => {
    console.error('Speech error:', e);
    setIsSpeaking(false);
  };
  window.speechSynthesis.speak(utterance);
};
```

### 3.5.2 Stop Functionality

Speech playback can be halted at any time using the `handleStop` function:

```
const handleStop = () => {
  window.speechSynthesis.cancel();
  setIsSpeaking(false);
};
```

### 3.5.3 Voice Selection Logic

The application prioritizes Google TTS voices when available, as they typically offer superior naturalness and clarity compared to system voices. A fallback chain ensures that the user always has a voice selected even in restricted environments:

1. Google US English voice (en-US, Google)
2. Google Hindi voice (hi-IN, Google)
3. Any available Hindi voice (hi-IN)
4. First available voice in the list

### 3.6 Database Integration

The frontend communicates with the backend API using Axios. The save operation posts the story title and content, and the fetch operation retrieves the list of saved stories:

```
// Save a story
const saveStory = async () => {
  if (!storyTitle || !extractedText) return;
  try {
    await axios.post('/api/save', {
      title: storyTitle,
      content: extractedText
    });
    fetchStories(); // Refresh the list
    setStoryTitle("");
  } catch (err) {
    console.error('Save error:', err);
  }
};

// Fetch all stories
const fetchStories = async () => {
  try {
    const { data } = await axios.get('/api/stories');
    setSavedStories(data.stories || []);
  } catch (err) {
    console.error('Fetch error:', err);
  }
};
```

```

}
};

```

## 4. TESTING AND EVALUATION

### 4.1 Testing Methodology

The ITIS Converter was evaluated using a multi-layer testing strategy encompassing unit testing, integration testing, performance testing, and user acceptance testing. This approach ensures that individual components behave correctly in isolation, components interact correctly when composed, the system performs within acceptable parameters under realistic conditions, and the final system meets user experience requirements.

Testing was conducted across three primary browsers — Google Chrome 120, Mozilla Firefox 120, and Microsoft Edge 119 — on Windows 11 and macOS Ventura. Mobile testing was performed on Chrome for Android and Safari for iOS. The test image corpus consisted of 50 images across five categories: printed text documents (10), handwritten notes (10), low-resolution scans (10), multilingual documents with Hindi and English (10), and natural scene text (10).

### 4.2 Unit Testing

#### 4.2.1 Backend API Tests

Backend route handlers were tested using Jest and Supertest, a library for testing Node.js HTTP servers. The following test cases were implemented:

| Test ID | Description                    | Input                             | Expected Output     | Result |
|---------|--------------------------------|-----------------------------------|---------------------|--------|
| UT-01   | POST /api/save with valid data | { title: 'T1', content: 'Hello' } | 201 + success: true | PASS   |
| UT-02   | POST /api/save missing title   | { content: 'Hello' }              | 400 + error message | PASS   |
| UT-03   | POST /api/save missing content | { title: 'T1' }                   | 400 + error message | PASS   |
| UT-04   | GET /api/stories returns array | None                              | 200 + stories array | PASS   |
| UT-05   | GET /api/stories empty DB      | None                              | 200 + empty array   | PASS   |
| UT-06   | Server starts on correct port  | None                              | Port 5000 listening | PASS   |

#### 4.2.2 Frontend Component Tests

React component tests were implemented using Jest and React Testing Library. Tests focused on verifying component rendering, user interaction handlers, and state update logic:

- ImageUpload: File type validation, file size limit enforcement, OCR trigger on valid file selection.
- SpeechControl: Rate slider produces correct utterance.rate values, pitch slider produces correct utterance.pitch values, play button calls speechSynthesis.speak().
- StoriesList: Renders correct number of story cards, clicking a story card populates the text area.

### 4.3 Integration Testing

Integration tests verified the end-to-end communication between the React frontend and the Express backend through the complete save and retrieve story workflow:

### 4.4 Performance Testing

#### 4.4.1 OCR Processing Time

OCR processing time was measured for images of varying resolution and content complexity. The following results represent the average of five runs per category on a mid-range laptop (Intel Core i5-11th Gen, 8GB RAM, Chrome 120):

| Test ID | Scenario                 | Steps   | Result |
|---------|--------------------------|---|--------|
| IT-01   | Full save workflow       | 1. Enter title 2. Enter text 3. Click Save 4. Verify in DB          | PASS   |
| IT-02   | Retrieve stories on load | 1. Add 3 stories to DB 2. Load app 3. Verify 3 cards shown          | PASS   |
| IT-03   | Backend unavailable      | 1. Stop backend 2. Click Save 3. Verify error handling              | PASS   |
| IT-04   | OCR then save            | 1. Upload image 2. Wait for OCR 3. Save result 4. Verify in stories | PASS   |

#### 4.2.2 API Response Time

Backend API response times were measured using Postman under simulated concurrent load:

| Endpoint         | Concurrent Users | Avg Response (ms) | P95 Response (ms) | Error Rate |
|------------------|------------------|-------------------|-------------------|------------|
| POST /api/save   | 1                | 45                | 67                | 0%         |
| POST /api/save   | 10               | 112               | 198               | 0%         |
| GET /api/stories | 1                | 23                | 41                | 0%         |
| GET /api/stories | 10               | 67                | 134               | 0%         |
| GET /api/stories | 50               | 210               | 445               | 0.20%      |

### 4.3 User Acceptance Testing

User Acceptance Testing (UAT) was conducted with a group of 15 participants consisting of 5 university students, 4 faculty members, 3 visually impaired individuals (using screen readers alongside the application), and 3 non-technical users. Participants were asked to complete three tasks: upload an image and verify the extracted text, adjust speech settings and play the text, and save a story and verify it appeared in the stories list.

| Criterion                   | Metric      | Score (out of 5) | Comments                       |
|-----------------------------|-------------|------------------|--------------------------------|
| Ease of image upload        | User rating | 4.3              | Drag-and-drop well received    |
| OCR accuracy (clear images) | User rating | 4.5              | Very accurate for printed text |
| OCR accuracy (handwriting)  | User rating | 3.1              | Struggled with cursive         |
| Speech naturalness          | User rating | 3.8              | Google voices preferred        |
| UI intuitiveness            | SUS Score   | 4.2              | Minimal learning curve         |
| Overall satisfaction        | User rating | 4.1              | Highly recommended             |

## 5. RESULTS AND DISCUSSION

### 5.1 OCR Accuracy Results

The OCR accuracy evaluation revealed a strong positive correlation between image quality and recognition accuracy, consistent with findings in the OCR research literature. Printed text on high-contrast backgrounds achieved accuracy rates exceeding 97%, demonstrating that Tesseract.js with its LSTM engine is suitable for production-grade text extraction from document images.

The overall Character Error Rate (CER) across all 50 test images was 8.6%, corresponding to a character-level accuracy of 91.4%. This is consistent with published benchmarks for Tesseract 4 LSTM mode on mixed image quality datasets. The Word Error Rate (WER) was 14.2%, reflecting the higher impact of recognition errors at word boundaries.

### 5.2 Speech Quality Evaluation

Speech quality was evaluated using a simplified Mean Opinion Score (MOS) methodology, where participants rated the naturalness and intelligibility of synthesized speech on a 1–5 scale. Google TTS voices achieved an average MOS of 3.9, while native system voices scored 3.2 on average. The speech rate and pitch controls were well-utilized by users, with 73% of participants adjusting the default speech rate — most commonly reducing it to 0.8x for improved comprehension.

### 5.3 System Performance Metrics

The system demonstrated acceptable performance characteristics for the target use case. Initial page load time (Largest Contentful Paint) measured 1.8 seconds on a standard broadband connection, within the recommended threshold of 2.5 seconds. The Tesseract.js WASM binary (approximately 10MB) is loaded asynchronously after the initial paint and cached by the browser's service worker or standard cache after the first load, reducing subsequent load overhead to under 200ms.

The MongoDB query performance for the stories retrieval endpoint averaged 23ms with an index on the date field, demonstrating that the current data model and indexing strategy are appropriate for the expected data volume (hundreds to low thousands of stories per user in a single-user deployment).

Memory usage during OCR processing peaked at approximately 180MB for a 4MP image on Chrome, which is within the acceptable range for modern desktop browsers. Mobile browsers showed higher memory pressure, occasionally triggering the garbage collector during processing of very large images, which caused intermittent UI freezes of 200–500ms duration.

## 6. CONCLUSION AND FUTURE SCOPE

### 6.1 Conclusion

This research paper has presented the design, implementation, and evaluation of the Image-to-Text-to-Speech Converter, a full-stack web application that demonstrates the practical integration of Optical Character Recognition and Text-to-Speech synthesis technologies in a modern browser environment. The project successfully achieved all seven primary objectives defined in Chapter 1, delivering a functional, tested, and documented accessibility tool.

The use of Tesseract.js as a WebAssembly-compiled OCR engine represents a significant architectural choice with important practical implications: by processing images entirely within the browser, the system eliminates the privacy concerns, latency overhead, and operational costs associated with server-side image processing. The Web Speech API provides a robust TTS layer that, while not matching the naturalness of premium cloud TTS services, is entirely free, widely supported, and sufficient for the accessibility use cases targeted by this application.

The evaluation results confirm that the system achieves an average character recognition accuracy of 91.4% across diverse image types, with accuracy exceeding 97% for high-quality printed document images. These results are consistent with published benchmarks for

Tesseract 4's LSTM engine and demonstrate that the system is suitable for practical accessibility use, particularly for users dealing with digitized printed documents, scanned reports, and photographed text.

The MongoDB-backed persistence layer provides reliable story management with sub-50ms API response times under typical load conditions. The React frontend offers an intuitive, responsive interface that received a System Usability Scale score equivalent of 4.2 out of 5 in user acceptance testing, indicating a high degree of usability.

The identified challenges — particularly WebAssembly loading performance, browser speech synthesis inconsistencies, and reduced accuracy for handwritten and scene text — represent well-understood limitations of the current implementation and have been mapped to specific enhancement proposals in Chapter 9. The proposed roadmap provides a clear path from the current academic prototype to a production-grade accessibility platform.

From a broader perspective, this project demonstrates that the modern web platform — with its WebAssembly runtime, Speech API, and powerful JavaScript frameworks — has matured to the point where complex AI-assisted accessibility tools can be delivered to users without requiring native application installation, cloud API dependencies, or specialized hardware. This democratization of accessibility technology aligns with the global imperative to reduce barriers to information access for all users.

In conclusion, the ITIS Converter achieves its stated goal of providing a practical, extensible, and well-documented example of full-stack web development applied to real-world accessibility needs, and provides a solid foundation for future enhancement into a production-grade assistive technology platform.

## 6.2 Future Scope

1. Image Preprocessing Pipeline
2. OCR Progress and Confidence Display
3. Error Handling and User Feedback
4. User Authentication and Multi-User Support
5. Audio Export Functionality
6. Multilingual OCR Support
7. Long-Term Enhancements (6–12 Months)
8. Cloud AI Integration
9. Progressive Web App (PWA) Conversion
10. Document Layout Analysis

## REFERENCES

1. Smith, R. (2007). An Overview of the Tesseract OCR Engine. In Proceedings of the 9th International Conference on Document Analysis and Recognition (ICDAR 2007), Vol. 2, pp. 629–633. IEEE.
2. LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278–2324.
3. Shi, B., Bai, X., & Yao, C. (2016). An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11), 2298–2304.
4. van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., ... & Kavukcuoglu, K. (2016). WaveNet: A Generative Model for Raw Audio. arXiv preprint arXiv:1609.03499.
5. Wang, Y., Skerry-Ryan, R., Stanton, D., Wu, Y., Weiss, R. J., Jaitly, N., ... & Saurous, R. A. (2017). Tacotron: Towards End-to-End Speech Synthesis. *Interspeech 2017*, 4006–4010.
6. Zen, H., Toda, T., Nakamura, M., & Tokuda, K. (2007). Details of the Nitech HMM-based speech synthesis system for the Blizzard Challenge 2005. *IEICE Transactions on Information and Systems*, 90(1), 325–333.
7. Breuel, T. M., Ul-Hasan, A., Al-Azawi, M. A., & Shafait, F. (2013). High-performance OCR for printed English and Fraktur using LSTM networks. In 2013 12th International Conference on Document Analysis and Recognition (pp. 683–687). IEEE.
8. Miesenberger, K., Karshmer, A., Penaz, P., & Zagler, W. (Eds.). (2012). *Computers Helping People with Special Needs. ICCHP 2012. Lecture Notes in Computer Science*, vol 7382. Springer, Berlin, Heidelberg.
9. Petrie, H., & Bevan, N. (2009). The evaluation of accessibility, usability and user experience. *The Universal Access Handbook*, 20, 1–16.
10. Reul, C., Christ, D., Hartelt, A., Balbach, N., Wehner, M., Springmann, U., ... & Puppe, F. (2019). OCR4all — An Open-Source Tool Providing a (Semi-)Automatic OCR Workflow for Historical Printings. *Applied Sciences*, 9(22), 4853.
11. World Health Organization (2023). Blindness and vision impairment. WHO Fact Sheet. Retrieved from <https://www.who.int/news-room/fact-sheets/detail/blindness-and-visual-impairment>

12. W3C Web Content Accessibility Guidelines (WCAG) 2.1. (2018). W3C Recommendation. Retrieved from <https://www.w3.org/TR/WCAG21/>
13. Mozilla Developer Network (2023). SpeechSynthesis — Web APIs. MDN Web Docs. Retrieved from <https://developer.mozilla.org/en-US/docs/Web/API/SpeechSynthesis>
14. Tesseract.js Documentation (2023). Tesseract.js: Pure JavaScript OCR. Retrieved from <https://tesseract.projectnaptha.com/>
15. Mongoose ODM Documentation (2023). Mongoose v7.x Documentation. Retrieved from <https://mongoosejs.com/docs/>
16. Facebook Inc. (2023). React 18 Documentation. Retrieved from <https://react.dev/>
17. OpenJS Foundation (2023). Express 4.x API Reference. Retrieved from <https://expressjs.com/en/4x/api.html>
18. MongoDB Inc. (2023). MongoDB Manual — CRUD Operations. Retrieved from <https://www.mongodb.com/docs/manual/crud/>
19. Twelve-Factor App Methodology (2017). The Twelve-Factor App. Retrieved from <https://12factor.net/>
20. Nielsen, J. (1994). Usability Engineering. Morgan Kaufmann. ISBN: 978-0125184069.