

# OCR検索可能PDF変換Webアプリ - 完全仕様書

バージョン 1.0.0

最終更新: 2026-1-15

リポジトリ: <https://github.com/11921604/OCR-PDF-Converter>

GitHub Pages (UIデモ): <https://11921604.github.io/OCR-PDF-Converter/>

関連ドキュメント:

- 仕様: <https://github.com/11921604/OCR-PDF-Converter/blob/main/docs/001-OCR-PDF-Converter/qsqs.md>
- 技術要件: <https://github.com/11921604/OCR-PDF-Converter/blob/main/docs/001-OCR-PDF-Converter/requirements.md>
- 計画: <https://github.com/11921604/OCR-PDF-Converter/blob/main/docs/001-OCR-PDF-Converter/plam.md>
- タスク: <https://github.com/11921604/OCR-PDF-Converter/blob/main/docs/001-OCR-PDF-Converter/tasks.md>

## 目次

- [概要](#)
- [システム構成](#)
- [機能仕様](#)
- [実装仕様](#)
- [API仕様](#)
- [データ仕様](#)
- [テスト方針](#)
- [デプロイ/配布](#)
- [運用/保守](#)

## 1. プロジェクト概要

### 1.1 背景と目的

スマートフォンでPDFファイルは、見出しは文書として認識できても、コンピュータにとって単純な画像の集合体です。このため、テキスト検索や複製ができます。情報の活用が制限されます。

#### 本アプリケーションの目的:

- スキャンPDFを検索可能なPDFに変換
- ブラウザ内で完結する処理 (プライバシー保護)
- 誰でも無料で利用できるWebアプリケーション

### 1.2 主要機能

- PDFアップロード:** ブラウザから直接PDFファイルをアップロード (50MB以下)
- OCR処理:** Pythonバックエンドで複数OCRエンジン (OmniOCR 2025.5、PaddleOCR 2.7.0.3) を並列実行し、最高精度の結果を自動選択
- 複数ページ対応:** 複数ページのPDFのバッチ処理とリアルタイム進捗表示
- 画像フォーマット対応:** JPEG、PNG、TIFFフォーマットを自動的にPDFに変換
- 検索可能PDF出力:** 検索とコピー可能なPDFファイルのダウンロード

### 1.3 技術スタック

レイヤー	技術	バージョン
フロントエンド	React	18.2.0
PDFレンダリング	pypdfium2 (Python)	4.30.0
OCRエンジン1	OmniOCR (Python)	2025.5
OCRエンジン2	PaddleOCR (Python)	2.7.0.3
PDF生成	pdf + ReportLab	5.1.0 + 4.2.0
バックエンド	Flask	3.0.0
ビルドツール	Webpack	5.104.1
テスト	Jest + Cypress	29.7.0 + 13.6.0
ホスティング	GitHub Pages (フロントのみ)	-
CI/CD	GitHub Actions	-

#### 1.3.1 OCRエンジンと並列処理アーキテクチャ

本システムは、複数OCRエンジンを並列実行し、最高精度の結果を自動選択する独自機構を実装:

##### 処理フロー:

- UIで複数エンジン選択 (OmniOCR、PaddleOCRのチェックボックス)
- 各PDFページで選択された全エンジンを並列OCR実行
- 各エンジンでの平均信頼度 (confidence) を計算
- 最も高い平均信頼度を保持エンジンの結果を採用
- 採用された結果をテキストレイヤーを生成

実装位置: [backend/main.py:4630-4672](#)

##### 統計情報:

- 各エンジンでの平均信頼度
- 総ページ数抽出
- 処理ページ数
- 最終エンジン名

この機構により、ドキュメントの種類や品質に応じて最適なOCRエンジンが自動選択されます。

#### 1.3.2 SSL証明書検証エラー対応

企業プロキシ環境やファイアウォール配下で、PaddleOCRがモデルファイルのダウンロード時にSSL証明書検証エラーが発生する場合があります。

##### エラーメッセージ例:

```
SSLERROR: [SSL: CERTIFICATE_VERIFY_FAILED] certificate verify failed: self signed certificate in certificate chain
```

実装済み対応 ([backend/main.py:4140](#)):

##### 1. SSL証明書検証の無効化

```
import ssl
ssl._create_default_https_context = ssl._create_unverified_context
```

##### 2. urllib警告の抑制

```
import urllib3
urllib3.disable_warnings(urllib3.exceptions.InsecureRequestWarning)
```

##### 3. 環境変数の設定

```
os.environ['REQUESTS_CA_BUNDLE'] = ''
os.environ['CURL_CA_BUNDLE'] = ''
```

##### 4. 複数エンジンチェック

```
ensure_paddleocr_available() 関数内でrequests.getに verify=False パラメータを自動注入
```

この対応により、自己署名証明書や企業内部CAを使用する環境でもPaddleOCRが正常に動作します。

### セキュリティと開発事項

この設定はローカル開発環境でのみ使用してください

プロダクション環境では適切なCA証明書を設定することと推奨します

モジュールダウンロード元 ([paddleocr/bcebos.com](#)) が信頼できることを確認してください

### 1.4 プロジェクト構造

```
OCR-PDF-Converter/
├── public/                  # 静的ファイル
│   ├── index.html          # HTMLテンプレート
│   ├── manifest.json       # PWAマニフェスト
│   └── sw/                  # サービスワーカー
├── components/             # Reactコンポーネント
│   ├── FileLoader.jsx      # ファイルアップロード
│   ├── OCRProgress.jsx     # 進捗表示
│   ├── DownloadButton.jsx  # ダウンロードボタン
│   └── App.jsx              # ビジネスロジック
├── services/                # サービス
│   ├── pdfProcessor.js     # PDF処理
│   ├── pdfGenerator.js     # カスタムReactフック
│   ├── hooks/              # useFileUpload, useOCR, js
│   └── utils/               # ユーティリティ関数
├── validators/              # fileValidator.js
├── coordinateConverter.js   # 座標変換
├── errorHandler.js         # エラーハンドラー
├── styles/                  # スタイルシート
├── index.css                # main.css
├── index.jsx                # エントリーポイント
├── tests/                   # テストファイル
│   ├── unit/               # ユニテスト
│   ├── integration/         # 統合テスト
│   └── e2e/                # E2Eテスト
├── docs/                   # ドキュメント
│   ├── README.md           # 本README
│   ├── 完全仕様書.md       # 完全仕様書
│   └── specs/               # 仕様書
│       ├── 001-OCR-PDF-Converter/
│       │   ├── spec.md
│       │   ├── requirements.md
│       │   ├── plan.md
│       │   ├── research.md
│       │   ├── data-model.md
│       │   ├── quickstart.md
│       └── tasks.md
├── .github/                 # GitHub Actionsワークフロー
│   ├── workflows/          # ワークフロー
│   ├── pages.yml           # パッケージ.json
│   ├── package.json        # package.json
│   ├── webpack.config.js   # webpack.config.js
│   └── README.md           # README
└──
```

## 2. 機能仕様

### 2.1 ユーザーストーリー

ユーザーストーリー 1: PDFアップロードとOCR処理 (P1)

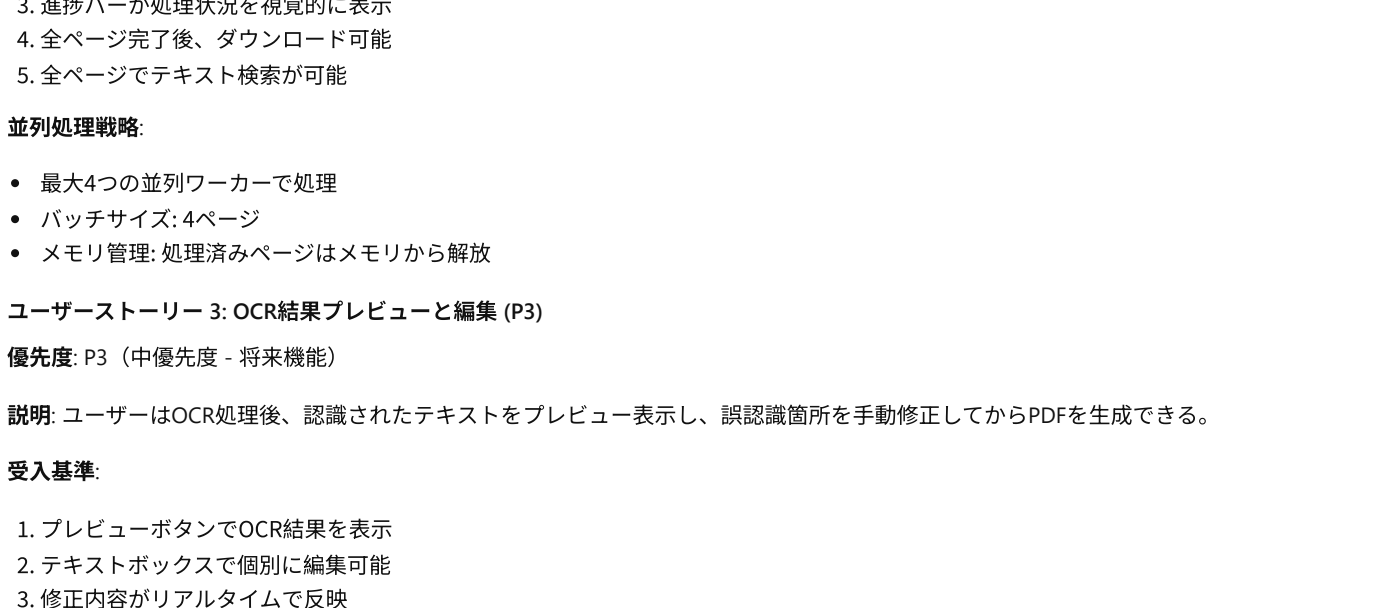
優先度: P1 (最高優先度 - MVP)

説明: ユーザーは1枚の複数ページを含むPDFファイルをブラウザからアップロードし、OCR処理を実行して、検索可能なテキストレイヤーを含むPDFをダウンロードできる。

#### 受入基準:

- PDFファイル (10MB以下) を選択できる
- ファイル名とファイルサイズが表示される
- OCR処理が完了するまで待機可能
- 処理完了後、ダウンロードボタンが表示される
- ダウンロード後、PDFでテキスト検索 (Ctrl+F) が可能

#### 処理フロー:



ユーザーストーリー 2: 複数ページPDFのバッチ処理 (P2)

優先度: P2 (最高優先度)

説明: ユーザーは複数ページを含むPDFファイルをアップロードし、全ページに対してOCR処理を一括実行できる。進捗状況がリアルタイムで表示される。

#### 受入基準:

- 複数ページPDF (10ページまで) を処理できる
- 「ページX/Y処理中...」形式で進捗を表示
- 進捗バーが処理状況を視覚的に表示
- 全ページ完了後、ダウンロードが可能
- 全ページでテキスト検索が可能

#### 並列処理戦略:

- 最大4つの並列ワーカーで処理
- バッチサイズ: 4ページ
- メモリ管理: 処理済みページはメモリから解放

ユーザーストーリー 3: OCR結果プレビューと編集 (P3)

優先度: P3 (中優先度)

説明: ユーザーはOCR結果後、認識されたテキストをプレビュー表示し、誤認識箇所を手動修正してからPDFを生成できる。

#### 受入基準:

- プレビューボタンでOCR結果を表示
- テキストボックスで個別に編集可能
- 修正内容がリアルタイムで反映
- PDF生成ボタンで最終PDF出力

ユーザーストーリー 4: 画像ファイル対応 (P2)

優先度: P2 (最高優先度)

説明: ユーザーはJPEG、PNG、TIFF形式の画像ファイルを直接アップロードでき、システムが自動的にPDFに変換してOCR処理を実行する。

#### 受入基準:

- JPEG、PNG、TIFF画像を選択できる
- 画像がPDFに自動変換される
- 変換後のPDFでOCR処理が実行される
- 画像サイズに応じて適切なPDFページサイズが設定される

ユーザーストーリー 5: 多様なページサイズ対応 (P2)

優先度: P2 (最高優先度)

説明: ユーザーはA4以外のページサイズ (A3、Letter、Legal、B4等) のPDFをアップロードでき、システムが自動的に適切な解像度でOCR処理を実行する。

#### 受入基準:

- A3、A4、B4、Letter、Legal、カスタムサイズに対応
- ページサイズを自動検出
- 各サイズに応じて300dpi基準で正規化
- すべてのサイズで正確なOCR処理が可能

## 2.2 エッジケースとエラーハンドリング

エッジケース	対応
10MBを超えるPDF	エラーメッセージ「ファイルサイズは10MB以下にしてください」を表示
非対応ファイル形式	エラーメッセージ「対応形式: PDF, JPEG, PNG, TIFFのみ」を表示
破損したPDF	エラーメッセージ「PDFファイルが破損しています」を表示
処理タイムアウト	10秒/ページでタイムアウト、再試行を促す
メモリ不足	エラーメッセージ「メモリ不足です。ページ数を減らしてください」を表示
既に検索可能なPDF	メッセージ「このPDFは既に検索可能です」を表示 (OCRスキップ)

## 3. 技術要件

### 3.1 機能要件

ID	要件	優先度
FR-001	PDF/画像ファイル (10MB以下) をブラウザからアップロード	P1
FR-002	PDFページを300dpi基準で画像としてレンダリング	P1
FR-003	OCR処理で日本語テキストの正確な認識を抽出	P1
FR-004	透明OCRテキストレイヤーを元のPDFに重ねて合成	P1
FR-005	検索可能なPDFのダウンロード	P1
FR-006	複数ページPDFの進捗表示	P2
FR-007	エラー通知の明示的な表示	P1
FR-008	完全クライアントサイド処理 (バックエンド不要)	P1
FR-009	ブラウザ内のみで処理 (サーバー送信なし)	P1
FR-010	ファイル形式検証とエラーメッセージ	P1
FR-011	多様なページサイズ対応 (A3/A4/Letter/Legal等)	P2
FR-012	画像ファイル (JPEG/PNG/TIFF) の自動PDF変換	P2

### 3.2 非機能要件

#### 3.2.1 パフォーマンス

メトリクス	目標値	測定方法
1ページOCR処理時間	5秒以内 (P95)	PerformanceObserver API
10ページOCR処理時間	50秒以内 (P95)	自動テスト
ページ読み込み時間	3秒以内 (P95)	Lighthouse
メモリ使用量 (ピーク)	2GB以下	Chrome DevTools
バンドルサイズ	5MB以下 (WASM含む)	webpack-bundle-analyzer

#### 3.2.2 信頼性

- 可用性: 99.9% (GitHub Pagesの可用性に依存)
- エラー率: 1%未満 (正常なPDFファイルで)
- OCR精度: 90%以上 (標準的な印刷品質)

#### 3.2.3 互換性

ブラウザ	最小バージョン	対応状況
Chrome	100+	完全対応
Firefox	100+	完全対応
Edge	100+	完全対応
Safari	15+	完全対応

#### 3.2.4 セキュリティ

- データ送信: 外部サーバーへの送信なし
- ブラウザ内処理: 全処理をクライアントサイドで完結
- HTTPS: GitHub Pagesで自動提供
- CSP: Content Security Policy設定済み
- データ削除: 処理完了後、メモリから自動削除

### 3.3 技術的制約

- ファイルサイズ: 最大10MB
- 対応画像: 非対応 (OCR)、日本語 (UI)
- 必須技術: WebAssembly対応が必要 (Service Workerでオフライン可能)

## 4. 実装計画

### 4.1 アーキテクチャ

#### 4.1.1 システムアーキテクチャ



#### 4.1.2 コンポーネント設計

##### プレゼンテーション層:

- App.jsx:** ルートコンポーネント、状態管理
- FileLoader.jsx:** ファイルアップロードUI
- OCRProgress.jsx:** 進捗バー表示
- DownloadButton.jsx:** ダウンロードボタン

##### ビジネスロジック層:

- pdfProcessor.js:** PDF読み込みと画像変換
- ocrEngine.js:** OCR処理の実行
- pdfGenerator.js:** テキストレイヤー生成とPDF合成

##### ユーティリティ層:

- fileValidator.js:** ファイル検証
- coordinateConverter.js:** 座標変換
- errorHandler.js:** エラーハンドリング

#### 4.2 実装フェーズ

##### Phase 1: Setup (T001-T014)

- プロジェクト構造作成
- 依存関係インストール
- 設定ファイル作成 (webpack, babel, ESLint等)

##### Phase 2: Foundational (T015-T032)

- ユーティリティ層実装
- サービス層実装
- コンポーネント作成

##### Phase 3: User Story 1 - MVP (T033-T050)

- Reactコンポーネント実装
- カスタムフック実装
- 統合テストとE2Eテスト

##### Phase 4: User Story 2 - Batch Processing (T051-T059)

- 複数ページ処理実装
- 進捗表示機能
- バッチ処理テスト

##### Phase 5: User Story 3 - Preview & Edit (T060-T070)

- プレビュー機能 (結果表示)
- 編集機能 (テキスト編集)

##### Phase 6: Image Support (T071-T080)

- JPEG/PNG/TIFF対応
- 画像→PDF変換
- 画像処理テスト

##### Phase 7: Page Sizes (T081-T088)

- A3/Letter/Legal対応
- ページサイズ自動検出
- 多様なサイズテスト

##### Phase 8: Polish (T089-T106)

- パフォーマンス最適化
- アクセシビリティ改善
- ドキュメント整備
- GitHub Pagesデプロイ

## 5. データモデル

### 5.1 主要エンティティ

#### 5.1.1 PDFFile

```
interface PDFFile {
  name: string; // ファイル名
  size: number; // サイズ (バイト)
  mimeType: string; // MIMEタイプ
  ArrayBuffer: ArrayBuffer; // バイナリデータ
  pageCount: number; // ページ数
}
```

#### 5.1.2 PDFPage

```
interface PDFPage {
  pageNumber: number; // ページ番号 (1始まり)
  width: number; // 幅 (ポイント)
  height: number; // 高さ (ポイント)
  canvas: HTMLCanvasElement; // レンダリング結果
  imageData: ImageData; // 画像データ
}
```

#### 5.1.3 OCRResult

```
interface OCRResult {
  items: OCRItem[]; // 認識されたテキスト項目
  confidence: number; // 全体の信頼度 (0-100)
  imageHeight: number; // 画像の高さ (座標変換用)
}
```

#### 5.1.4 OCRItem

```
interface OCRItem {
  text: string; // テキスト内容
  x: number; // 幅 (ポイント)
  y: number; // 高さ (ポイント)
  width: number; // 幅 (PDFポイント)
  height: number; // 高さ (PDFポイント)
  fontSize: number; // フォントサイズ
  confidence: number; // 信頼度
}
```

#### 5.2 データフロー



## 6. テスト戦略

### 6.1 テストピラミッド



#### 6.2 テストカバレッジ目標

レイヤー	カバレッジ目標
ユニテスト	80%以上
統合テスト	70%以上
E2Eテスト	主要フロー100%

#### 6.3 テストツール

- ユニテスト: Jest + @testing-library/react
- 統合テスト: React Testing Library
- E2Eテスト: Cypress
- パフォーマンステスト: Lighthouse CI

#### 6.4 主要テストケース

##### ユニテスト

- fileValidator: MIME type検証、ファイルサイズ検証
- coordinateConverter: 座標変換、フォントサイズ計算
- errorHandler: カスタムエラークラス、エラーメッセージ変換
- pdfProcessor: PDF読み込み、ページレンダリング
- ocrEngine: OCR実行、結果ファイルタギング
- pdfGenerator: テキストレイヤー生成、PDF合成

##### 統合テスト

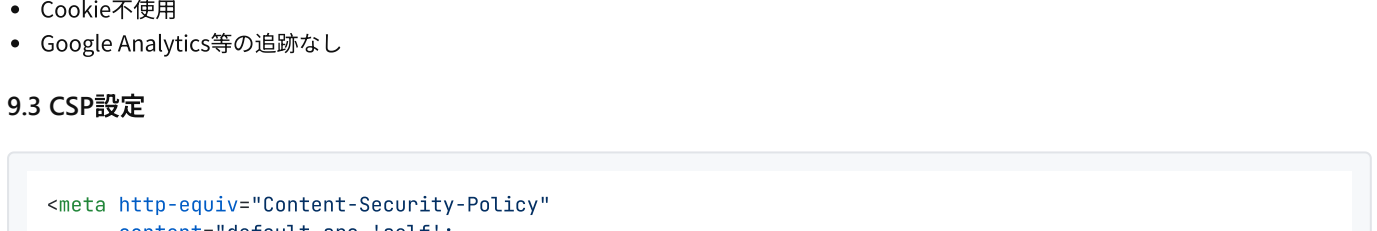
- PDFアップロード → OCR → ダウンロード
- 複数ページPDF処理
- 画像ファイル処理
- エラーハンドリング

##### E2Eテスト

- 正常系: PDFアップロード → OCR処理 → ダウンロード
- エラー系: 非対応ファイル、サイズ超過
- レスポンス: モバイル、タブレット、デスクトップ
- アクセシビリティ: キーボード操作、ARIA属性

## 7. デプロイメント

### 7.1 デプロイアーキテクチャ



#### 7.2 GitHub Actions ワークフロー

ファイル: [.github/workflows/pages.yml](#)

##### トリガー:

- push:** ブランチへのプッシュ
- manual:** workflow\_dispatch

##### ジョブ:

- build:**
  - Node.js 18セットアップ
  - 依存関係インストール
  - テスト実行
  - 100%コードカバレッジ取得
  - アーティファクトアップロード
- deploy:**
  - GitHub Pagesへデプロイ

#### 7.3 デプロイURL

本環境: <https://11921604.github.io/OCR-PDF-Converter/>

#### 7.4 環境変数

変数名	値	用途
PUBLIC_URL	/OCR-PDF-Converter	GitHub Pagesのサブパス
NODE_ENV	production	ビルドモード

## 8. パフォーマンスとスケーラビリティ

### 8.1 パフォーマンス最適化戦略

#### 8.1.1 バンドル最適化

- Code Splitting:** Reactコンポーネントの遅延読み込み
- Tree Shaking:** 未使用コードの削除
- Minification:** JS/CSSの圧縮
- Gzip Compression:** GitHub Pagesで自動有効化

#### 8.1.2 OCR処理最適化

- Web Workers:** OCR処理をメインスレッドから分離
- 並列処理:** 最大4ページを並列処理
- メモリ管理:** 処理済みページを即座に解放
- 進捗的レンダリング:** ページ単位で結果を表示

#### 8.1.3 キャッシング戦略

- Service Worker:** アプリケーションのオフライン動作
- モデルキャッシュ: OCRエンジンのモデルファイルをキャッシュ (.paddle, .onnx)
- CDN: PDF.jsをCDNから読み込み

### 8.2 スケーラビリティ

- ユーザーあたり1プロセス (ブラウザ制御)
- 100%コードカバレッジ取得
- GitHub Pages: 無制限トラフィック (リミット: 100GB/月)

#### 将来のスケーリング:

- Progressive Web App (PWA) 化
- IndexedDBでの処理履歴保存
- Web WorkerのDynamic Scaling

## 9. セキュリティとプライバシー

### 9.1 セキュリティ原則

- データの外部送信禁止:** 全処理をブラウザ内で完結
- HTTPS通信:** GitHub Pagesで自動提供
- Content Security Policy:** XSS攻撃対策
- 依存関係の脆弱性管理:** npm auditで定期チェック

### 9.2 プライバシー保護

#### データ処理方針:

- アップ