

防衛省防衛研究所仕様書			
件名	戦史史料のOCR学習データ等の作成	作成	戦史研究センター

## 1. 適用範囲

この仕様書は、防衛省防衛研究所が所蔵する戦史史料の利活用に資する学習用データ等の作成について規定する。

## 2. 役務に関する要求

### 2. 1 作業内容

契約相手方は以下の作業を実施する。

同所が保持する旧陸海軍文書の画像データを作業側に貸与する。(官側は当該画像データを格納したハードディスクを準備、貸与する。)

#### (1) 基礎史料調査

防衛研究所戦史研究センター史料室が保有する史料について、現地(当研究所F1棟史料庫)での目視による調査を含め、史料の概要を調査する(詳細は第3項に示す。)

#### (2) AI-OCR用学習データ作成

契約相手側は、貸与された画像から様々な書式パターンの適切な画像を抽出し、手書き文字計50万字以上のAI(Artificial Intelligence)-OCR(Optical Character Recognition)学習用データセットを作成する。(詳細は第4項に示す。)

#### (3) AI-OCR精度向上実証調査

前項で作成した学習データをAI-OCRシステムに学習させ、当該システムを開発し、OCRの精度向上効果について報告書としてまとめる。(詳細は第5項に示す。)

#### (4) 活字識別精度向上策

NDL-OCRを使用した際の活字識別精度向上に資する方策を講じ、この際、追加学習データを作成する場合は、2.1(2)項に含めて提出する。(詳細は第6項に示す。)

#### (5) AI-OCRによるテキスト化

約2,000点の手書き文字資料及び活字史料に対して、AI-OCRを用いてテキスト化処理を実施する。(詳細は第7項に示す。)

#### (6) 生成AIを活用した検索・分析システムの実証環境の提供

契約相手方は、官側がテキスト化された史料の検索や内容の整理・分析を可能にする、生成AIを活用したシステムを利活用できることを確認するための環境を準備する。(詳細は第8項に示す。)

### 2. 2 作業期間

契約締結日から令和8年3月23日(月)まで

### 2. 3 作業場所

契約相手方作業所等

ただし、官側から要請があった場合は速やかに協議する必要があるため、作業場所は関東甲信越地方に限定する。

## 2. 4 準 拠

契約相手方は、契約後速やかに官側と作業打合せを行い、作業内容等の詳細について確認決定し、作業スケジュール表を作成し、官側の承認を得て作業を実施する。

作業スケジュール表は、作業項目ごとに作業人員、想定処理量、作業期間等を記述した週単位のスケジュールを記載し、速やかに官側に提出すること。

## 2. 5 権 利

本作業で作成された基礎史料調査の報告書、学習用データ、精度向上実証の報告書、活字史料のOCRテキストの編集著作権は官側に帰属する。また学習用データは官側が許諾した場合のみ、AI-OCRシステムに学習させることができるものとする。

## 2. 6 秘密保持

本作業を実施するにあたり知り得た情報及び作成したデータ、報告書等の内容については、第三者に洩してはならない。

## 2. 7 個人情報保護

本作業を実施するにあたり、官側が貸与する画像データ内に個人情報を確認した場合、個人情報保護法や関係法令を遵守し、個人情報データの漏洩防止のために必要な安全管理措置を講じつつ、当該史料は官側に返却するものとする。

## 2. 8 成果物

契約相手方は以下に示す成果物を官側に納品するものとする。格納媒体は契約相手方で準備するものとする。細部については、官側との協議に基づき、実施するものとする。

- (1) 基礎史料調査報告書：紙及びデータ格納媒体 [DVD-R]
- (2) 学習用データ※：データ格納媒体 [DVD-R] ※
- (3) 精度向上実証調査報告書：データ格納媒体 [DVD-R]
- (4) AI-OCRテキスト：データ格納媒体 [DVD-R]
- (5) 生成AIを活用した検索・分析システムの利用手順書  
※2. 1 (4) により生じたデータを含む。

## 2. 9 品質保証

- (1) 官側へ納品した成果物に不具合が発見された場合は、納入後3年間、契約相手方の責任において速やかに修正等の対応を行うものとする。
- (2) 修正作業等により電子画像データの引渡しを受ける場合は官側の許可を得るものとする。

## 3 基礎史料調査

### 3. 1 史料現地調査

官側との協議に基づき、官側が保有する史料全点（約75,000点、約2,500万枚）について調査を行い、必要に応じて当研究所F1棟史料庫において現

地調査を行う。史料の規模、性質、レイアウト、画像の画質、筆跡の難易度などについて調査を行い、AI-OCRテキスト化作業を行うにあたって、テキスト化、史料の整理、分類等を含め必要な作業工程や時程などについて報告書にまとめる。またAI-OCRテキスト化の作業を行う上で効率的なグループに仕分ける

### 3. 2 手書き史料・活字史料の仕分け

官側が保有する史料の画像データにつき、手書き史料と活字史料を自動的に仕分けるシステムまたはプログラムについて検討・製作する。そのシステムまたはプログラムにより、史料約2,000点(100万枚)について仕分けを行う。対象画像は3.1項の調査に基づき官側と協議の上、選定するものとする。その際、個人情報が含まれる画像は選外とする。

## 4 AI-OCR学習データ作成

4.1 基礎史料調査において自動振り分けされ、分類された手書き史料のうち、官側との協議によって先行的に学習データを作成するものについて、崩し字(行書及び草書体)を概ね6割以上含むAI-OCR学習データを50万文字分以上作成する。その際、個人情報が含まれる画像は選外とする。

4.2 AI-OCR学習データセットの仕様を以下に示す。

- (1) 仕様に未記載の事項で疑問が生じた場合は、官側と協議の上決定するものとする。
- (2) 評価用データは、画像データと画像データに対応するJSONデータの組み合わせで作成する。
- (3) 画像データが0.5°以上傾いている場合は傾きを補正し、文字が極力正立した状態に修正する。画像データの傾き補正情報は下記の構造で作成し、TSV形式テキスト(またはExcelデータ)で納品する。

[TSV形式テキスト出力例]

画像名 1. jpg[タブ] 0.6[改行]

画像名 2. jpg[タブ]-1.2[改行]

画像名 3. jpg[タブ]-1.2[改行]

...

画像名 600. jpg[タブ]2.0[改行]

- (4) JSONデータには画像に対応したアノテーション情報(座標情報、テキスト等)を入力する。
- (5) 学習用データは画像のレイアウト種別等のグループごとにそれぞれ1フォルダ(以下:グループフォルダ)に格納し、以下の構造で作成する。

```
グループフォルダ└─「img」フォルダ└─画像名 1. jpg (傾き補正後の画像)
                    │
                    │   └─画像名 2. jpg
                    │   └─  ...
                    │   └─画像名 600. jpg
                    │
                    └─「json」フォルダ└─画像名 1. json
                                        └─画像名 2. json
```

ト …  
└─画像名 600. json

画像補正データ.xlsx

- (6) JSON データのテキストは官側が別途指示する翻刻方針を基準として入力する。翻刻方針の変更が必要な場合は、官側と協議の上決定するものとする。
- (7) JSON データのテキストは入力精度 99%以上とする。  
ただし、史料の破損や汚れ等で判読不能な文字、及び翻刻方針上入力不能な文字については、官側と協議の上、精度保証の対象外とする。
- (8) JSON データは、a「資料情報 (BOOK)」、b「画像情報 (PAGE)」、c「行矩形情報 (LINE)」の3階層で構成する。
- a BOOK の定義
- BOOK は「資料名 (TITLE)」の属性を持つ。
  - BOOK は子要素として PAGE を持つ。
- b PAGE の定義
- PAGE は「対象画像ファイル名 (NAME)」「高さ (H)」「横幅 (W)」の各属性を持つ。
  - PAGE は子要素として LINE を持つ。
- c LINE の定義
- LINE は内包する文字の集合の外接矩形とする。
  - LINE は「行種別 (TYPE)」「書字方向 (ALLIGN)」「行内テキスト (STR)」「4 点座標 (POINTS)」の各属性を持つ。
  - TYPE は、“本文”・“ルビ”のいずれかを入力する。
  - ALLIGN は、“縦組”・“横組 (左から右)”・“横組 (右から左)”・“反転”のいずれかを入力する。
  - STR は、行内の文字を OCR または目視で解読したテキストを入力する。テキストは校正することで精度 99%以上を保証する。
  - STR に史料の破損や汚れ等で判読不能な文字、及び翻刻方針上入力不能な文字が含まれる場合には当該文字を「=」で入力する。
  - POINTS は「x1 y1 x2 y2 x3 y3 x4 y4」の形式で入力する。

#### [JSON 構造例]

```
BOOK {資料名, [  
  PAGE {幅, 高さ, 画像ファイル名, [  
    LINE {座標情報, テキスト, 付加情報※},  
    LINE {座標情報, テキスト, 付加情報},  
    LINE {座標情報, テキスト, 付加情報}  
  ]}  
]}
```

]]

※本文/ルビの区分、書字方向情報 (縦組、横組 (左から右)、横組 (右から左)、反転)

#### [JSON 形式 OCR テキスト出力例]

```
{"BOOK":  
  {"TITLE": "海軍省-公文備考-S1-43-3396", "PAGE": [  
    {"W": "3348", "H": "4708", "NAME": "KS31_0013_01.TIF", "LINE": [  
      {
```

```

        {"POINTS":[20, 340, 940, 340, 940, 1200, 20, 1200], "ALLIGN":"縦組",
"TYPE":"ルビ", "STR":"デンキキョウカイ"},
        {"POINTS":[2139, 361, 2139, 551, 2269, 551, 2269, 361], "ALLIGN":"縦組", "TYPE":"本文", "STR":"社団法人電気協会主催ノ下ニ"},
        {"POINTS":[1380, 372, 1380, 714, 1540, 714, 1540, 372], "ALLIGN":"縦組", "TYPE":"本文", "STR":"制轉機ノ轉把ヲ廻ハシテ...緩ム"}
    ]}
  ]}
}

```

## 5 AI-OCR精度向上実証調査

第4項で作成したAI-OCR学習データをAI-OCRに学習させ、学習後のAI-OCRの処理結果の精度をF値(※)及び編集距離で評価する。評価用データについては官側が提供するものとする。使用するAI-OCRプログラム及び実行環境は契約相手方で準備する。また、AI-OCRプログラムは近代の草書を含む手書き文字に対して100万字以上学習済みであり、第4項で作成した学習データを追加で学習させた場合の精度向上に関して報告するものとする。

※ 国立国会図書館によるF値(Fmeasure)の定義は以下の通り。

```

ytrue   = {正解文字情報に含まれる文字の多重集合},
Ypred   = {認識結果に含まれる文字の多重集合}
Precision = |ypred ∩ ytrue| / |ypred|, Recall = |ypred ∩ ytrue| / |ytrue|
Fmeasure = 2 * Recall * Precision / (Recall + Precision)

```

また、精度以外の定性的な情報(レイアウト認識の影響等)についても分析する。評価・分析結果はAI-OCR精度向上実証調査報告書にまとめて納品する。

## 6 NDL-OCRプログラムの改修

官側が保有する活字史料のうち、主要なレイアウトを抽出し、レイアウト認識の向上等を目的としたNDL-OCRプログラムの改修を行う。レイアウトの抽出のための史料の点数及び種類については官側と協議するものとする。プログラムの改修を行った後、官側が提供する評価データに基づき精度評価を行い、その結果を官側に報告する。その内容については、「AI-OCR精度向上実証調査報告書」に含むものとする。その際、認識精度については中央値90パーセント以上を目標とするが、その結果については官側と協議するものとする。また改修にあたって、追加学習データを作成する場合、第4項作業の成果物と併せてこれを提出する。

## 7 AI-OCRによるテキスト化

第3項の自動仕分けによって抽出された手書き文字史料及び活字史料に対して、それぞれ8万枚以上の画像を選定し、AI-OCRプログラムによるテキスト化を行う。その際、手書き文字史料に対しては第5項で検証したAI-OCRを使用し、活字史料に対しては第6項で改修したNDL-OCRを使用する。

## 8 生成AIを活用した検索・分析システムの実証環境の提供

クラウド上でセキュリティを担保した生成AIを活用し、第7項のテキストデータ及

び官側が別途支給したテキストデータによる、全文検索等利活用の検証を行うための環境を構築、これの利用手順書を作成し提供する。生成AI及びそれを利用するサービスについては、官公庁及び地方自治体への導入実績のあるサービスを使用するものとする。

## 9 成果物、提出書類の納期等

### 9.1 成果物

成果物引き渡しは表1のとおり。

表1

名 称	数 量	引渡時期	引渡場所
基礎史料調査報告書	1 式	作成完了次第	防衛省 防衛研究所
AI-OCR学習データ	1 式		
AI-OCR精度向上実証調査報告書	1 式		東京都新宿 区市谷本村 町5-1
AI-OCRテキストデータ	1 式		
生成AIを活用した全文検索システムの利用手順書	1 式		

### 9.2 その他提出書類

提出書類については表2のとおりとする。

表2

名 称	部数	提出時期	様式等	提出場所	備 考
従業員名簿及び資格証明書	1 部	契約後速やかに	様式随意及び写し	防衛省防衛研究所	本業務に携わる作業員名簿及び作業実績
作業スケジュール表	1 部		様式随意	東京都新宿区市谷本村町5-1	本業務の作業予定

## 10 入札資格要件

入札参加希望者の資格要件は以下のとおりとし、要件を満たすことを証明する資料を官側の指定する日までに書面で提出し、承認を得るものとする。

### 10.1 実 績

公文書館または公文書館等に類する機関が保有する古文書や公文書の電子化・テキスト化業務に関し、過去5年間でAI-OCRの学習・評価用データとして活字1000万字程度又は手書き100万字程度の役務計3回以上の作成実績及び当該機関向けのAI-OCRの開発・改修実績を有すること。

### 10.2 精度評価レポート

官側が提供する史料2点（手書き史料1点、活字史料1点）に付き、国会図書館が公開しているNDL-OCR及びNDL-古典籍OCRを使用し精度評価を行い、その結果を報告書として提出すること。

### 10.3 学習データ作成

官側が提供する手書き史料1点に付き、今回の契約において使用する予定のAI-OCRを使用して学習データを作成し、当該AI-OCRの精

度評価の結果と合わせて提出すること。

10. 4 セキュリティ

官側が貸し出す画像の保管場所については入退出管理及び24時間警備体制が確立され、使用パソコンについては、社外への流出防止処置が施されているものとする。

10. 5 作業従事者

精度評価の報告書作成は、実務経験5年以上の正社員が行うこと。また、当該者は精度評価データ及び学習データ作成の知識を有している、(社)日本文書情報マネジメント協会が行う上級文書情報管理士の資格保有者であること。

現場責任者及び現場副責任者は、上記実務において、責任者又は副責任者等を務めた実績を有する者であること。

11 その他

この仕様書に疑義が生じた場合は、速やかに官側と協議するものとする。