

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（約 2,000 万枚）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) 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 200 万文字以上の 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 3 項に示す。）。

(2) 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 4 項に示す。）

(3) AI-OCRによるテキスト化

官側が提供する史料のうち、2 万点以上、もしくは 100 万枚以上を選定し、それらを手書き・活字・その他に正確に分類し、そのうち手書きのものについてのみ、(2) で精度を向上させた AI-OCR によってテキスト化を行う。（詳細は第 5 項に示す。）

(4) 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 6 項に示す。）

2. 2 作業期間

契約締結日から令和 9 年 3 月 26 日（金）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、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、その旨を官側に通知し、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) AI-OCR学習用データ※ : データ格納媒体 [DVD-R]
- (2) AI-OCR精度向上実証調査報告書: データ格納媒体 [DVD-R]
- (3) AI-OCRテキストデータ: データ格納媒体 [DVD-R]
- (4) 生成AIを活用した検索・分析システムの利用手順書

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1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. AI-OCR学習データ作成

3. 1 官側との協議によって学習データを作成するものについて、活字を含まず、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200万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。

3. 2 AI-OCR用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) 学習用データは、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
                                     └─ ...
                                     └─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 (ALIGN)」「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。
 - ALIGN 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], "ALIGN": "縦  
組", "TYPE": "ルビ", "STR": "デンキキョウカイ"},  
      {"POINTS": [2139, 361, 2139, 551, 2269, 551, 2269, 361],  
"ALIGN": "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],  
"ALIGN": "縦組", "TYPE": "本文", "STR": "制轉機ノ轉把ヲ廻ハシテ...緩ム"}  
    ]}  
  ]}  
}
```

4 AI-OCR 精度向上実証調査

第3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手書き文字に対して100万字以上、(内草書50万字以上)を学習済みであり、第3項およびこれまで官側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 = $\frac{|y_{pred} \cap y_{true}|}{|y_{pred}|}$,

Recall = $\frac{|y_{pred} \cap y_{true}|}{|y_{true}|}$

Fmeasure = $2 \frac{Recall * Precision}{Recall + Precision}$

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。

5 AI-OCRによるテキスト化

官側が提供する史料のうち、2万点以上または100万枚以上の史料を官側と協議の上選定し、手書き／活字／その他（図、細かい表組等）に正確に分類整理する。

そのうち、手書き史料（最大100万枚）を選別し、第3項で作成した学習データで精度を向上させたAI-OCRを用いてテキスト化すること。テキスト化する史料の選別にあたっては、官側の承認を得ること。

テキスト化は官側が提供する「翻刻方針」に従うこと。その他、テキスト化作業の詳細については、官側と協議の上、その指示に従うこと。

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

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5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、検索・分析システムの実証環境を構築、この利用手順書を作成し提供する。実証環境の提供期間は、契約締結後できるだけ早期の時期から、令和9年3月26日までとする。

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。使用するサービスはISMAP登録済みの高度なセキュリティを担保したクラウド環境（日本リージョン限定）を用いたものとする。また以下の要件を満たすものとする。

- a. ハルシネーション対策：①事実が確認できない場合、正直に「分からない」と回答する。②推測や予測が含まれる場合、必ずその旨を明記する。③回答の根拠となったデータのソースを明示（引用）する機能を備えること。
- b. データガバナンス：入力データがAIモデルの学習に利用されないことを担保する技術構成をとること。
- c. 検索性の高度化：全文検索に加え、史料種別、作成年代等の属性情報を用いたフィルタリング機能を備えること。これにより、特定の属性に絞り込んだ上でのAIによる要約・分析を可能とすること。
- d. 実証に使うデータ量：本実証環境に取り込むデータ量は、戦史叢書（約7万ページ）、R7年度事業にてAI-OCRによるテキスト化したもの（約15万ページ）に加えて、第5項にてテキスト化するデータ（約100万ページ）をベースとすること。
- f. エージェント機能：単一の検索にとどまらず、生成AIが自律的に分析計画の立案、複数史料の比較検証、検索結果の妥当性評価、および情報不足時の自動的な再検索（反復検索）を繰り返すエージェント機能を備え、複雑な分析要求に対しても論理的な回答を導き出せること。

7 成果物、提出書類の納期等

7.1 成果物

納品する成果物は表1のとおり。

表 1

名 称	数 量	引渡時期	引渡場所
A I - O C R 用 学 習 データ	1 式	令和 9 年 3 月 2 6 日	防衛省 防衛研究所
A I - O C R 精 度 向 上 実 証 調 査 報 告 書	1 式		
A I - O C R テ キ ス ト データ	1 式		
生 成 A I を 活 用 し た 検 索 ・ 分 析 検 索 シ ス テ ム の 利 用 手 順 書	1 式		東 京 都 新 宿 区 市 谷 本 村 町 5 - 1

7. 2 その他提出書類

提出書類については表 2 のとおり。

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

8 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。

8. 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 A I - O C R の学習・評価用データとして近代（明治、大正、昭和）の手書き文字 1 0 0 万字以上の作成業務計 3 回以上の実績及び上記機関向けの A I - O C R の開発・改修実績を有すること。

8. 2 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 A I - O C R を使用して学習データを作成し、当該 A I - O C R の精度評価の結果と合わせて提出すること。

8. 3 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 2 4 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。

8. 4 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人員をあてること。

近世および近代の崩し字を含む手書き文字の解読業務について、経験年数 5 年以上の人員を入力作業従事者および品質管理者として計 5 名以上配置すること。

近世・近代の手書き文字等の難読文字の入力品質の確保のため、国文学または歴史学の専門家を監修者として配置すること。

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。

業務従事予定者の名簿を入札時に提出することとし、契約後に業務従事者の入れ替えを行う場合は、官側の承認を得るものとする。

9 その他

9. 1 本事業以前に行われた過去の成果物において、現作業に影響する不具合事項を新たに認めた場合、官側との協議によりその修正を検討するものとする。
9. 2 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。