

令和3年度 OCR処理プログラム研究開発作業 報告資料

資料作成：株式会社モルフォAIソリューションズ

事業の概要

本事業の目的・実施概要

本事業の目的

国会図書館デジタルアーカイブ画像のOCR
処理を活用したテキストデータの作成による、
本文検索サービスや視覚障害者等への
デジタル化資料の提供の利便性向上



実施概要

貴館の所蔵するデジタル化資料の画像を利用して、所定の性能を
満たすOCR処理プログラムの研究開発として以下二点を実施

① OCR学習用データセット及びレイアウト情報データセット作成

- 1870年代～1960年代 理系図書・文系図書
 - 1870年代～1990年代 雑誌
- 計17,375件を作成

② OCR処理プログラムの開発

- 本文テキスト領域の文字の認識精度向上及び、
画像のテキスト化処理速度の向上

本事業のアプローチ

■課題

貴館の所蔵資料は膨大であるため、開発するOCR処理プログラムには高い認識精度と高速な処理速度の両立が求められる。一方で従来手法には実用面の課題がある。例えば商用の日本語OCR処理プログラムは現代の資料を対象としているため、近代書籍に特有の複雑なレイアウトや文字配置、字形、旧字旧仮名遣いに十分に対応しておらず精度面で課題がある。また、近年、OCR処理において深層学習等の先端的な機械学習技術を取り入れた事例が注目を高めている。ただし、これらの多くは研究段階や試行段階のものであるため、精度が高い場合であっても処理速度が十分でない等の傾向がみられる。

上記を踏まえ、当社及び当社再委託先は本事業の目的に適い実用化に資するOCR処理プログラム実現のために以下のアプローチを採用する。



① OCR学習用データセット及びレイアウト情報データセットの作成

良質な学習データの構築

精度の高いOCR処理プログラム構築のためには良質な学習用データセットの構築が不可欠である。学習用データセットの構築を担う凸版印刷は、明治期以降の近代書籍に関するOCR処理の実績、及び字形データセット構築の実績を豊富に有しています。そのため過去の知見を活かし、貴館が保有する書籍画像データの特徴を踏まえ、複雑なレイアウトや特殊な文字配置、字形といった近代書籍の特徴を網羅した良質なデータセットの構築を目指す。

② OCR処理プログラムの開発

精度面と処理速度の面を両立した機械学習アルゴリズムの構築

一般に深層学習をはじめとする機械学習のアルゴリズムは、GPU等の豊富な計算リソースを前提に処理速度よりも精度面での工夫を行うものが多い傾向にある。一方で、モルフォ及びモルフォAIソリューションズでは、これまでスマートフォンや車載カメラといった計算資源の限られた端末での機械学習・画像処理プログラムの開発を多数行ってきた。したがって、モルフォが保有する技術・ノウハウを活かすことにより精度のみならず処理速度でも実用に足るOCR処理プログラムの開発を目指す。

実施内容と成果

OCR学習用データセットの実施概要

元画像データ抽出・選別 画像角度補正

実施内容

- ・ 貴館貸与画像データ約250万資料の集計
- ・ データセット用の画像抽出・選別（約250万資料→2.5万件）作業
- ・ 抽出・選別した画像の補正（複製/拡張子変更/傾き補正/トリミング等）

要件定義

- ・ 納品用データセットの詳細仕様の検討
（データ構成/文字種/行矩形要素/ブロック要素/インライン要素/欧文/数式・化学式/出力形式/矩形情報の並び順）

OCR学習用データセット構築

- ・ 要件定義で固めたOCR学習データセット仕様書に基づき、納品用のデータセットを作成
- ・ 7/29初回納品、以降9/27までに計6回、17,375件を納品（**納品済みのOCR学習用データセットについて修正対応含む**）
- ・ レイアウト情報データセットは、12/21初回納品、以降22/2/7までに計4回、27,878件を納品

アウト プット



- ・ OCR学習用データセット画像抽出手順書
- ・ OCR学習用データセット作成用画像_画像補正手順書



- ・ データセット構築用画像（2.5万件）



- ・ OCR学習用データセット作成仕様書



- ・ 納品データー式（アノテーションデータ/xmlデータ /（画像）補正情報

```
<?xml version="1.0" encoding="UTF-8" standalone="true"?>
<OCRDATASET xmlns="ND:OCRDATASET">
  <PAGE KYOKAKU="true" HEIGHT="3272" WIDTH="2255" IMAGENAME="R0000068_contents_L.jpg">
    <BLOCK HEIGHT="1296" WIDTH="20" Y="807" X="2062" TYPE="JL2"/>
    <LINE HEIGHT="2110" WIDTH="45" Y="744" X="2018" TYPE="本文" STRING="右五尺間半に面打出線分。且">
      <CHAR HEIGHT="39" WIDTH="40" Y="744" X="2019" MOJI="右"/>
      <CHAR HEIGHT="39" WIDTH="41" Y="799" X="2018" MOJI="左"/>
      <CHAR HEIGHT="40" WIDTH="38" Y="849" X="2020" MOJI="尺"/>
      <CHAR HEIGHT="36" WIDTH="32" Y="906" X="2023" MOJI="間">
      <CHAR HEIGHT="40" WIDTH="30" Y="957" X="2020" MOJI="半"/>
    </LINE>
  </PAGE>
</OCRDATASET>
```

補正情報、1960-9 - マモ帳

ファイル名	画像(E)	書式(O)	表示(O)	ヘルプ(H)
1334854_R0000072_contents_L.jpg	4370	3205	-0.8	2052
1334854_R0000072_contents_R.jpg	4324	3143	0	2252
1335051_R0000112_contents_L.jpg	5530	3794	0	2614
1335051_R0000112_contents_R.jpg	5584	3870	-0.78	2814
1335920_R0000076_contents_L.jpg	4338	3101	0	2105

OCR処理プログラムの実施概要

OCR処理プログラムの実施フロー



実施項目

- OCR処理プログラムに実装する手法の検討・検証
 - 分割
 - 傾き補正
 - レイアウト/行認識
- (上記の検討・検証を踏まえ) 実装手法の課題整理・対応
- プログラム設計 (内部/外部設計) ・開発※

中間報告

- 精度・速度改善
- 見出し認識・著者名対応

※ プログラム開発は10月以降も継続

中間報告・進捗報告

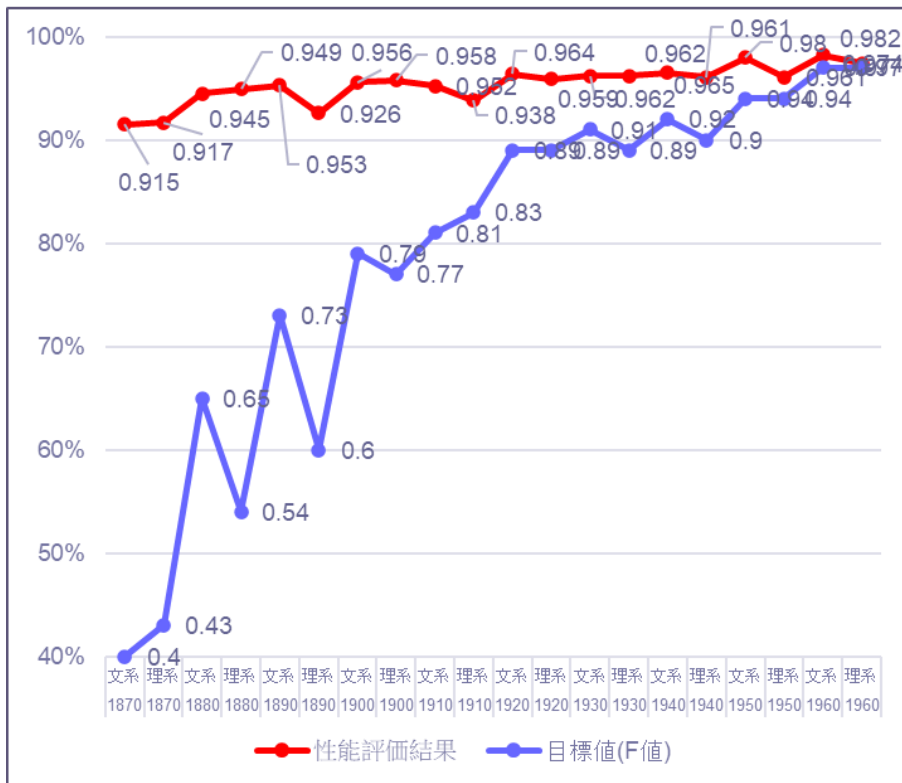
OCR処理プログラム認識性能の達成状況

- 全納品データ17,375件のOCR処理プログラムに実装する全手法（①見開き分割、②傾き補正、③レイアウト認識、④行認識）を統合した際の認識性能は以下の通り、30個（目標値達成）/30個（評価対象、参考値3個は除く）となり、仕様書上の認識精度基準は達成することができた。

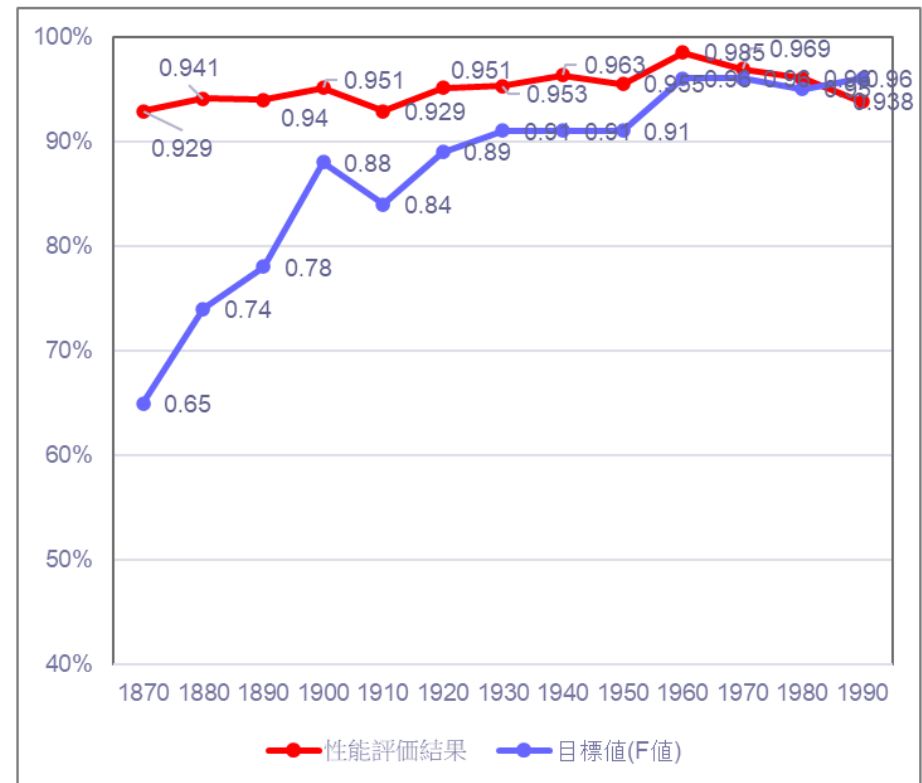
※各年代ごとに10件のスコアを算出した中の中央値で判定（全330件）

認識精度基準と認識性能評価結果まとめ

書籍種別・年代別の精度評価



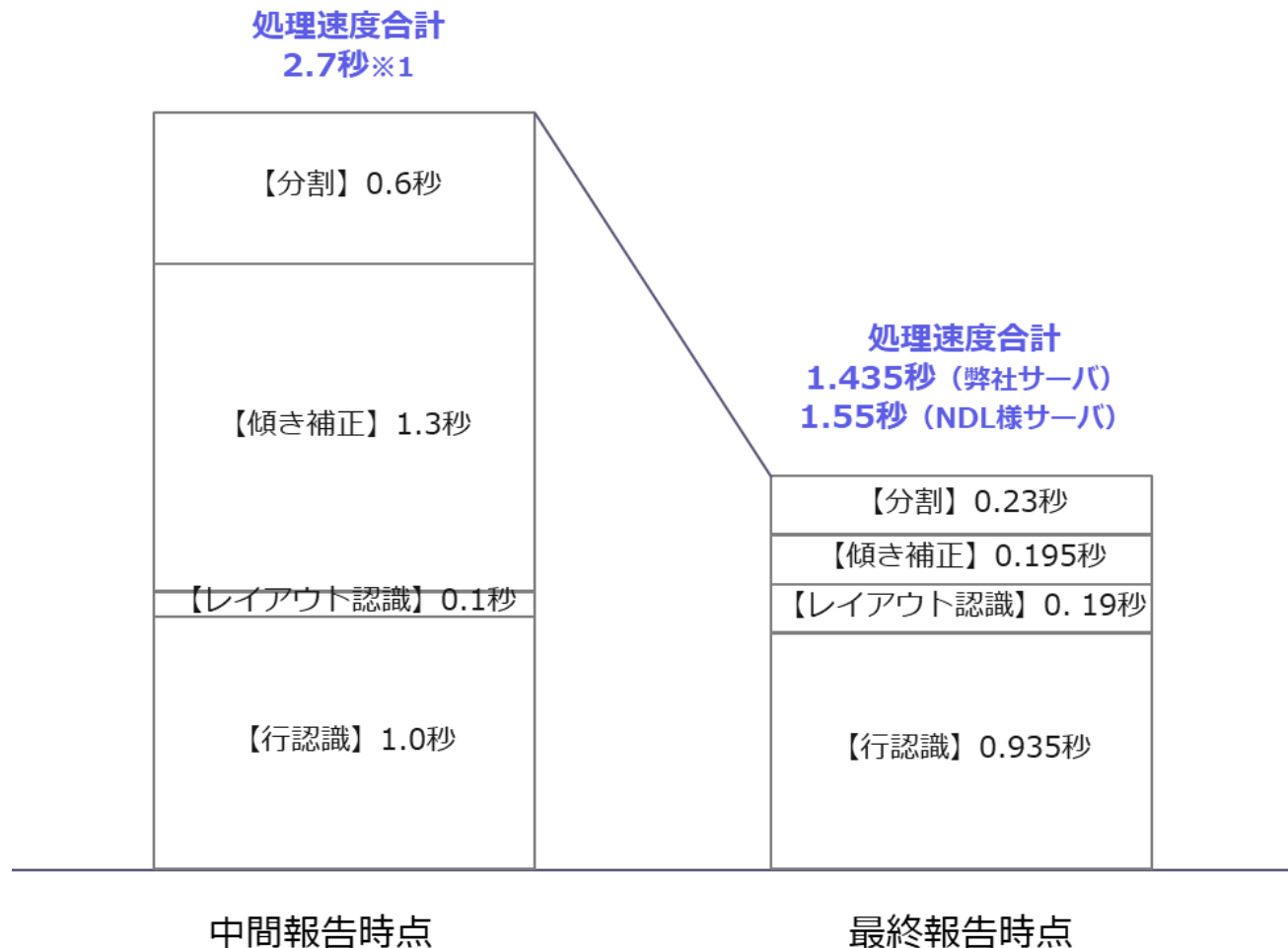
雑誌の年代別の精度評価
(1970-1990年代は参考値)



OCR処理プログラム速度性能の達成状況

- 中間報告時点から更に43%処理時間を削減し約1.55秒/枚の処理時間を達成、目標の2秒/枚を大きく上回る速度性能を達成

モジュール別・1枚あたり処理速度



※1 貴館と弊社で動作環境で異なること、分割と傾きは入出力を含めた値

OCR学習用データセット

※データセットの作成作業は凸版印刷株式会社が担当

【ご参考】OCR学習用データセット納品物のサンプル（1910年代図書）

imgファイル

本文テキスト内外の座標情報が記載されたxmlファイル

政治部八十四 下編 雑税中

五〇

```
<?xml version="1.0" encoding="UTF-8" standalone="true"?>
<OCRDATASET xmlns="NDLOCRDATASET">
  <PAGE KYOKAKU="true" HEIGHT="3272" WIDTH="2255" IMAGENAME="R0000068_contents_L.jpg">
    <BLOCK HEIGHT="1296" WIDTH="20" Y="807" X="2062" TYPE="ルビ"/>
    <LINE HEIGHT="2110" WIDTH="45" Y="744" X="2018" TYPE="本文" STRING="右五尺間半に而打出候分長横平均五尺壹間に付御年貢百文に定 壹間不詰尺者 六寸より貳尺" DIRECTION="縦">
      <CHAR HEIGHT="39" WIDTH="40" Y="744" X="2018" MOJI="右"/>
      <CHAR HEIGHT="35" WIDTH="41" Y="799" X="2018" MOJI="五"/>
      <CHAR HEIGHT="40" WIDTH="38" Y="849" X="2020" MOJI="尺"/>
      <CHAR HEIGHT="36" WIDTH="32" Y="906" X="2023" MOJI="間"/>
      <CHAR HEIGHT="40" WIDTH="39" Y="957" X="2020" MOJI="半"/>
      <CHAR HEIGHT="31" WIDTH="29" Y="1015" X="2025" MOJI="に"/>
      <CHAR HEIGHT="37" WIDTH="37" Y="1066" X="2021" MOJI="而"/>
      <CHAR HEIGHT="37" WIDTH="39" Y="1119" X="2022" MOJI="打"/>
      <CHAR HEIGHT="35" WIDTH="34" Y="1171" X="2025" MOJI="出"/>
      <CHAR HEIGHT="38" WIDTH="39" Y="1223" X="2022" MOJI="候"/>
      <CHAR HEIGHT="37" WIDTH="41" Y="1277" X="2022" MOJI="分"/>
      <CHAR HEIGHT="9" WIDTH="8" Y="1318" X="2053" MOJI="長"/>
      <CHAR HEIGHT="39" WIDTH="38" Y="1331" X="2023" MOJI="横"/>
      <CHAR HEIGHT="40" WIDTH="40" Y="1382" X="2023" MOJI="平"/>
      <CHAR HEIGHT="34" WIDTH="37" Y="1439" X="2022" MOJI="均"/>
      <CHAR HEIGHT="37" WIDTH="38" Y="1490" X="2022" MOJI="五"/>
      <CHAR HEIGHT="35" WIDTH="40" Y="1543" X="2022" MOJI="尺"/>
      <CHAR HEIGHT="38" WIDTH="40" Y="1596" X="2022" MOJI="壹"/>
      <CHAR HEIGHT="9" WIDTH="9" Y="1636" X="2049" MOJI="間"/>
      <CHAR HEIGHT="40" WIDTH="37" Y="1647" X="2024" MOJI="に"/>
      <CHAR HEIGHT="38" WIDTH="34" Y="1701" X="2025" MOJI="付"/>
      <CHAR HEIGHT="30" WIDTH="29" Y="1757" X="2026" MOJI="御"/>
      <CHAR HEIGHT="38" WIDTH="36" Y="1806" X="2024" MOJI="年"/>
      <CHAR HEIGHT="34" WIDTH="39" Y="1863" X="2023" MOJI="貢"/>
      <CHAR HEIGHT="39" WIDTH="39" Y="1915" X="2022" MOJI="百"/>
      <CHAR HEIGHT="40" WIDTH="40" Y="1966" X="2022" MOJI="文"/>
      <CHAR HEIGHT="39" WIDTH="37" Y="2019" X="2025" MOJI="に"/>
      <CHAR HEIGHT="38" WIDTH="40" Y="2072" X="2022" MOJI="定"/>
      <CHAR HEIGHT="31" WIDTH="28" Y="2129" X="2029" MOJI="壹"/>
      <CHAR HEIGHT="36" WIDTH="37" Y="2180" X="2023" MOJI="間"/>
      <CHAR HEIGHT="9" WIDTH="9" Y="2219" X="2051" MOJI="不"/>
    </LINE>
  </PAGE>
</OCRDATASET>
```

納品対象画像の補正情報

補正情報_1910-0 - メモ帳

ファイル(F)	編集(E)	書式(O)	表示(V)	ヘルプ(H)					
PID	画像名	回転	補正後横幅	補正後高さ	補正角度	トリミング矩形X座標			
1766303	R0000068_contents_L.jpg	5850	3996	4004	-0.19	2863			
1766303	R0000068_contents_R.jpg	5856	4004	3996	-0.28	3063			
1905974	R0000065_contents_L.jpg	5116	3786	3780	-0.21	2387			
1905974	R0000065_contents_R.jpg	5110	3780	3786	0.13	2587			
897690	R0000068_contents_L.jpg	4350	3272	3282	0.11	2055			
897690	R0000068_contents_R.jpg	4358	3282	3272	0.23	2255			

画像選別基準

画像選別については、貴館の内容確認の上、以下の抽出及び除外対象とした。

[抽出対象画像]

- ・ テキストと図版入り画像
- ・ テキストのみ画像
- ・ 図版のみ画像※キャプション付き
- ・ 表組入り画像
- ・ 全面表組
- ・ 版本

[除外画像]

- ・ 表紙、扉
- ・ 画質の悪い画像（解像度、極端な傾き、ピンボケ、折り目、汚れ、虫食い、ノドの開き不足）
- ・ 図版のみの画像※キャプションなし
- ・ 全面広告
- ・ グラビア（＋グラビア内テキスト）画像
- ・ ブランクページ
- ・ 片観音、両観音、Z折り、特殊折り
- ・ 全面外国語
- ・ 楽譜
- ・ マンガの吹き出し（雑誌・図書内のマンガ記事を含む）
- ・ 新聞の縮刷版
- ・ 目次

OCR学習用データセットの納品実績

貴館と合意した分割単位で均等に画像を抽出、企画提案時に考慮していなかった理系/文系が不明な図書は、全データを考慮して330資料を抽出するという方針に基づき、各年代の納品数量バランスを取りつつ、最終的に17,375件の納品を9/27までに完了させた。

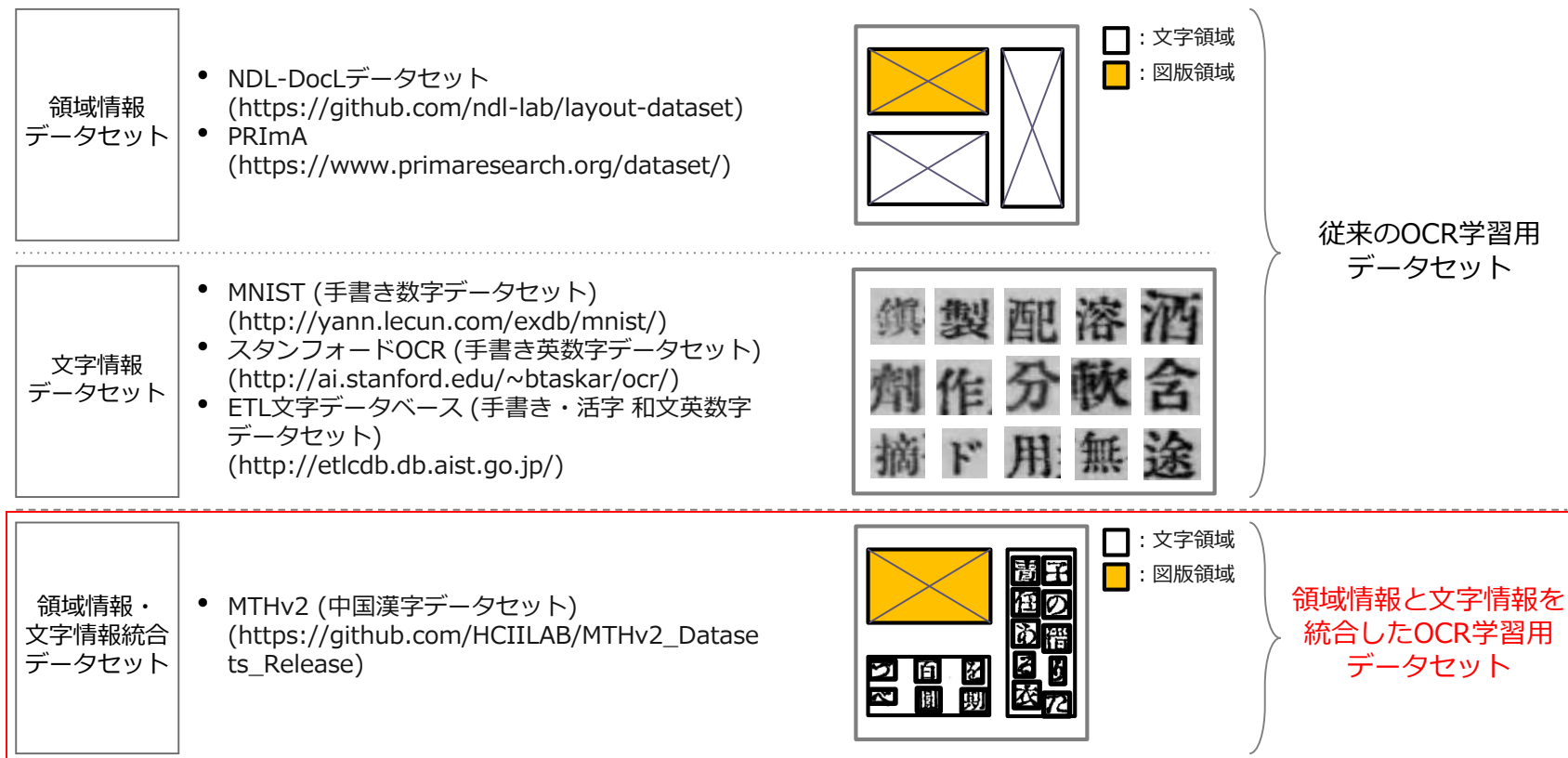
OCR学習用データセットの納品数量及び内訳（全体サマリー）

No	年代	7/29納品	7/30納品	8/16納品	8/30納品	9/13納品	9/27納品	合計
		実績	実績	実績	実績	実績	実績	実績
1	新字(1960年以降の理系図書・文系図書・雑誌)	521	614	1,007	1,173	913	1,570	5,798
2	旧字① (1910～1959年の理系図書・文系図書・雑誌)	534	520	1,112	971	786	2,166	6,089
3	旧字② (1909年以前の理系図書・文系図書・雑誌)	516	546	926	1,047	1,314	1,139	5,488
合計		1,571	1,680	3,045	3,191	3,013	4,875	17,375

OCR学習用データセット構成（1/2）

本事業におけるOCR学習用データセットは、様々なAI-OCRの実装に対応するため、領域に関する情報と文字に関する情報はまとめて1データでOCR学習用データセットを作成する。

■領域情報データセット、文字情報データセット、および両者を統合したデータセットとしては、それぞれ以下の例があげられる。



本事業でのOCR学習用データセット構成

OCR学習用データセット構成 (2/2)

OCR学習用データセットに含まれる情報の一覧を示す。本文テキスト内の行矩形情報、文字矩形情報に加えて、本文外の画像情報、基本的な文字種以外のブロック要素（図版、欧文、数式、化学式）についても以下の通り定める。

OCR学習用データセットに含まれる情報

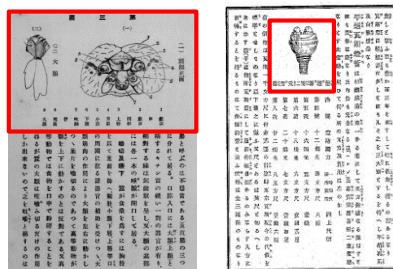
本文テキスト内	(本文内) 行矩形情報	<ul style="list-style-type: none"> 座標 (X,Y) 高さ (H) 横幅 (W) 書字方向 (縦、横、右から左) 行内文字列情報 (テキスト)
	(本文内) 文字矩形情報	<ul style="list-style-type: none"> 座標 (X,Y) 高さ (H) 横幅 (W) 文字情報 (文字コード)
本文テキスト外	画像情報	<ul style="list-style-type: none"> 画像名 画像高さ (H) 画像横幅 (W) 匡郭有無 傾き修正情報 (回転方向・角度) ※ トリミング情報 (左右・切断位置) ※ その他
	文字以外の領域情報	<ul style="list-style-type: none"> 座標 (X,Y) 高さ (H) 横幅 (W) 種類 (図版、欧文、数式、化学式など)

文字以外の領域情報

- 図版、数式、化学式は、本事業の目的である本文検索用テキスト作成においてテキスト認識の障害要因であり、精度向上のためには文字認識の対象外として排除することが望ましい。本件OCR学習用データセットでは、これらの要素については文字以外の領域情報として作成し、将来的にOCR処理プログラムの機能として排除処理を実装する際の学習用データとしての活用を想定する。
- 欧文は領域情報として、ブロック要素およびインライン要素でOCR学習用データセットに文字以外の領域情報として作成する。

※ 欧文は既存OCRソフトの認識精度が高いため、本事業の目的を鑑みて、本件OCR処理プログラムでは欧文領域の認識のみ実行し、別途既存OCRソフトで文字認識を実施する、という実装を想定する。
 欧文のOCR学習用データセットは既存のオープンデータが存在するため、本提案では新規に作成せず、既存オープンデータの使用を想定する。

[図版] (ブロック要素)



[欧文] (ブロック要素・インライン要素)

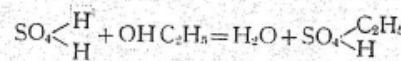
May 3, 1948 (commemoration day of the enforcement of the Japanese Constitution) shall be made a holiday for all government offices. However, the superintendent official of an office may at his own discretion, make officials under his charge attend to office work so that the management of business of urgent necessity may not be hindered.

等を一括して攻究する方寧ろ便利なるに
 然有機化学 (Organic chemistry) なる一部門を
 機物を無機物より分別して説述すると常

[数式] (ブロック要素・インライン要素)

$$\frac{1}{A} = \frac{1}{15} - \left(-\frac{1}{13} \right) = \frac{1}{7}$$

[化学式] (ブロック要素・インライン要素)

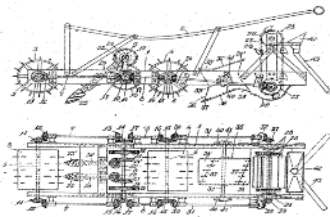


※画像補正作業で作成するTSVファイルに記載

OCR学習用データセット レイアウト情報—ブロック要素

- BLOCKの種類 (TYPE) は「図版」「表組」「柱」「ノンブル」「ルビ」「系統図」「数式」「化学式」「欧文ブロック」

図版



表組

表1 チタン半成品寸法範囲の一例

寸法 (mm)	重量 (kg)	半成品寸法 (mm)
直径	3.5	スラブ 厚さ: 100~260 幅: 600~1310
高さ	19.0	角ビレット 114~180角
		丸ビレット φ120~φ370

柱

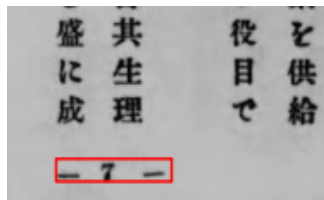


数式

$$\frac{1}{A} = \frac{1}{15} - \left(-\frac{1}{13} \right) = \frac{1}{7}$$

BLOCK

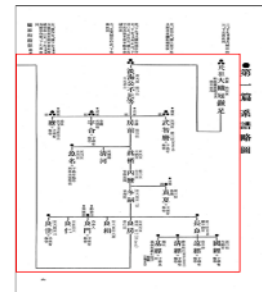
ノンブル



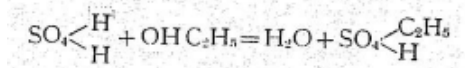
ルビ



系統図


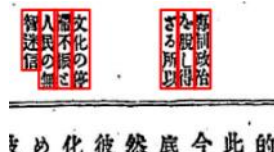


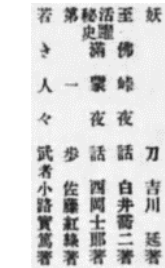

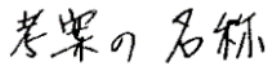
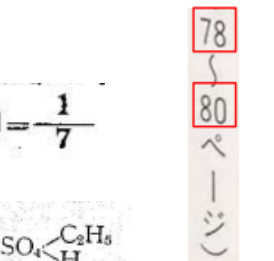
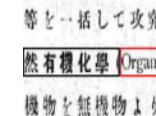
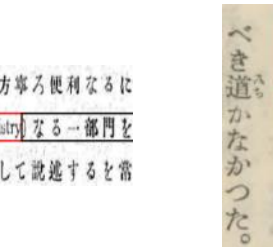
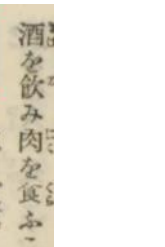
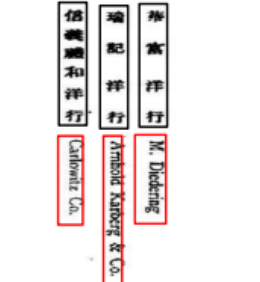


化学式



OCR学習用データセット レイアウト情報ーライン要素・インライン要素

- LINEの種類 (TYPE) は、「キャプション」「頭注」「割注」「見出し」「著者名」「本文」
- INLINEの種類 (TYPE)、「手書き」「縦中横」「数式」「化学式」「欧文」「色付き文字」「回転欧文」

LINE	キャプション	頭注	割注	見出し	著者名	本文
						
INLINE	手書き	縦中横	欧文	色付き文字	色付き文字	回転欧文
						

OCR学習用データセットの文字種

基本的な 文字種

- ・ ひらがな
- ・ カタカナ
- ・ 数字
- ・ JIS第一水準漢字・JIS第二水準漢字
- ・ 下記の記号
 - 半角記号 , . - / ()
 - 句読点 、 。
 - 括弧 () [] 『 』 【 】 「 」
 - 丸付き文字 ①②③④⑤⑥⑦⑧⑨⑩⑪⑫⑬⑭
 ㊦㊧㊨㊩㊪㊫㊬㊭㊮㊯
 - 丸付き漢字 ㊰
 - 括弧付き文字 (株)(名)(資)(有)
 - 繰り返し記号 \ / め ん 々 " > ズ ム ム
 - 図形 ○ ● □ ◆ ▲ △ ▽
 - その他の記号 — ・ ※ ↓ → ↑ ← ⇄ ⇅ ⇆ ⇇ ? ~ = ≠ …

凸版印刷による事前調査で出現率3,000位までの
JIS第一水準漢字・JIS第二水準漢字に包摂が可能な
JIS第二水準外の漢字

欧文・ ギリシア 文字

- ・ 本文中に出現する3文字程度までの欧文・ギリシア文字
- 半角アルファベット52文字 (U+0041-U+005A、
U+0061-U+007A)
- ギリシア文字48文字 (U+0391-U+03A9、U+03B1-
U+03C9)

留意事項

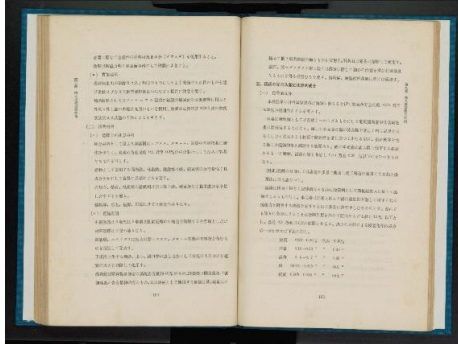
- ・ 欧文・ギリシア文字は、十分な領域がありレイアウト認識ができるものについては、矩形情報のみのOCR学習用データを作成し、文字情報は入力対象外とする。
- ・ 半角文字が存在する(ASCIIコード内)の英数字記号は、別紙2「性能評価対象とする資料の内訳及び判定基準」の仕様を踏まえ、すべて半角のコードポイントを割り当てて。
- ・ 合略仮名や結合文字、特殊な丸付文字など、コードポイントがないため1文字で入力不可能な字形は本件OCR学習用データ作成の対象外とする。
- ・ 訓点、図版内文字、表組内文字は文字種に関わらずOCR学習用データ作成の対象外とする。
- ・ 変体仮名は対応する現代仮名遣いのコードポイントを割り当てる。
漢字と変体仮名どちらとも取れる文字は漢字のコードポイントを割り当てる。
- ・ 行頭字下げ、文中スペースは学習用データ作成の対象外とする。
- ・ 包摂不可能なJIS第二水準外文字、および汚れやかすれなどで文字が判読不能な文字は「二」を割り当てる。

OCR処理プログラム

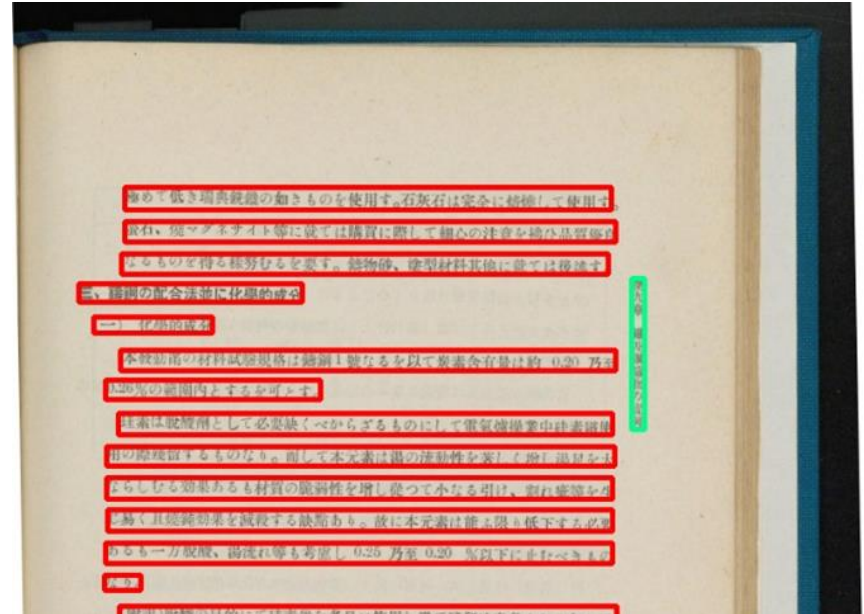
実装した手法全体の処理フロー概要

実装した手法全体の処理フローは以下の通りである。

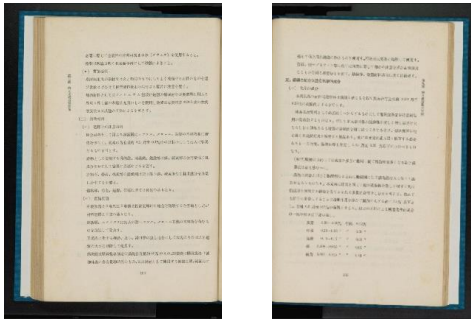
入力画像



③レイアウト認識(レイアウト認識モジュール)



①見開き分割(分割モジュール)



④行認識(行認識モジュール)

極めて低き瑞典銑鐵の如きものを使用す。石灰石は



極めて低き瑞典銑鐵の如きものを使用す。石灰石は

...

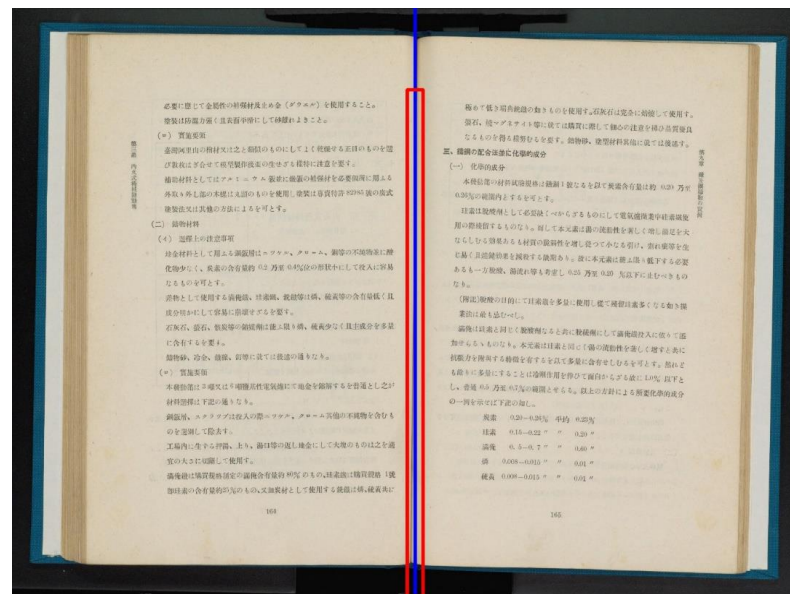
実装した各手法について 中間報告までの振り返り

アプローチ

- NDLラボで開発されたノド元分割プログラムを採用 (https://github.com/ndl-lab/ssd_keras)
- SSDを用いたノド元のx座標検出器
- 貴館よりモデルを提供いただき利用

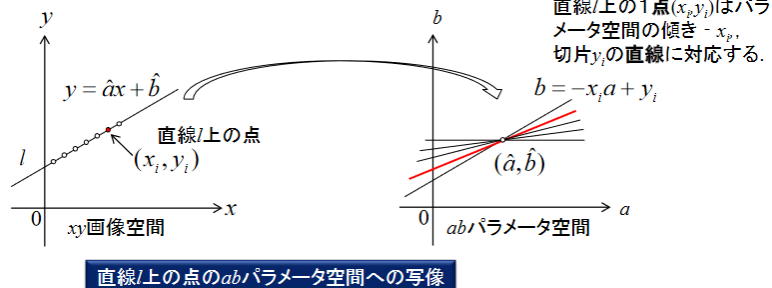
見開き分割

手法のイメージ

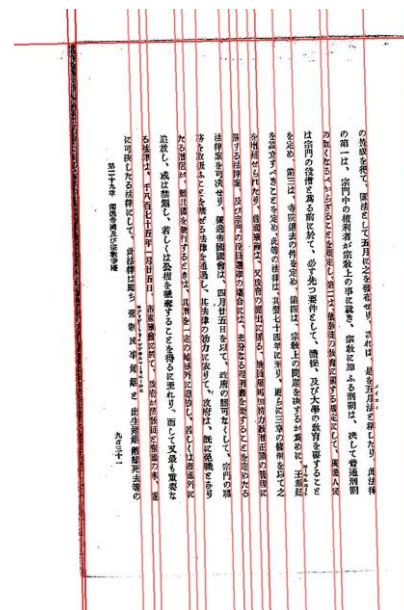


- Hough変換を用いた画像内の直線要素の検出
- 検出された直線の角度のうち最頻値をページの傾きとして採用
- Hough変換のイメージ

傾き補正



<http://www.cfme.chiba-u.jp/~haneishi/class/digitalgazo/6HoughTransform.pdf>



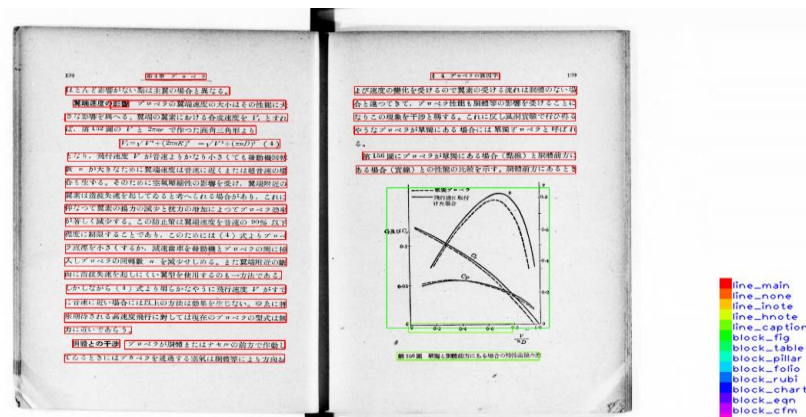
実装した各手法について 中間報告までの振り返り

アプローチ

手法のイメージ

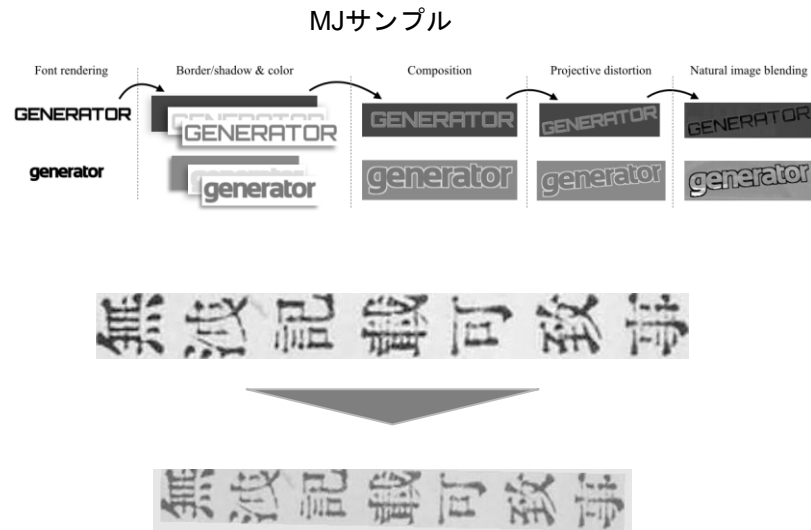
レイアウト 認識

- Cascade Mask RCNNのDetection部分を採用
以下のアプローチを参照
 - インドの古文書行認識 (Indiscapes)
(<https://arxiv.org/abs/1912.07025>)
 - 2020年SIGNATEコンペ (3rd, 4th)
(<https://signate.jp/competitions/218>)
- レイアウト認識特有のチューニング
 - 画像に対して対象が小さい → 入力解像度を大きく
 - 対象の数が多い → Region Proposal(途中段階の候補領域提示)の上限増
 - アスペクト比が極端(縦長、横長) → アンカーボックスのサイズ調整



行認識

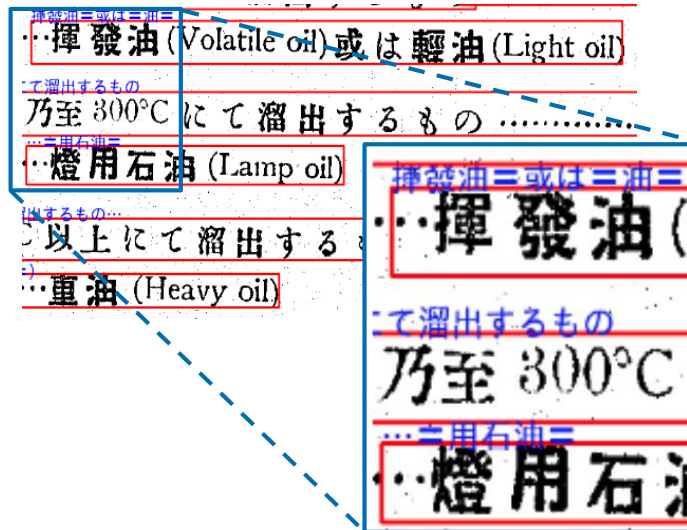
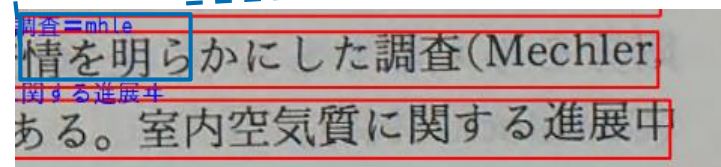
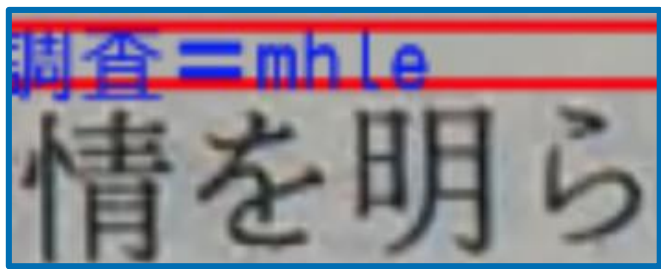
- シーンテキスト認識手法Rosettaを採用
(<https://arxiv.org/abs/1910.05085>)
以下のシーンテキスト認識手法の比較研究を参照
(<https://arxiv.org/abs/1904.01906>)
- backboneにResNet, decoderにCTCを利用
- 合成画像英単語データセットMJを利用
- 学習時の各種Augmentation
 - アスペクト比変更
 - パディング追加
 - 角度ずれ



中間報告時点での課題

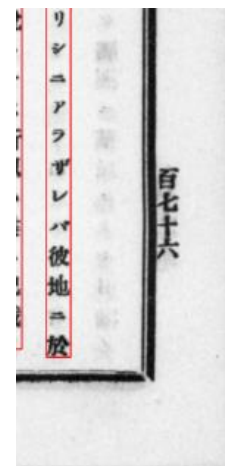
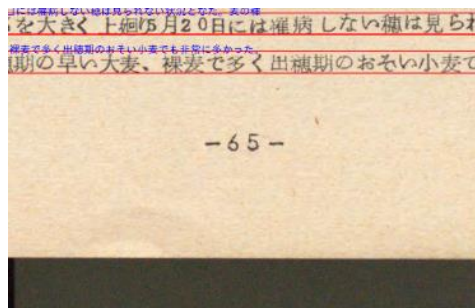
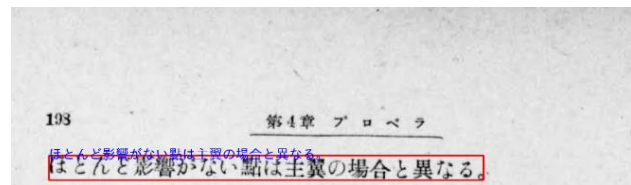
1. 特に理系文書における欧文、数値、記号などの認識性能改善

英文字、数値の認識欠落のサンプル



2. レイアウト、文字間違いの認識性能改善

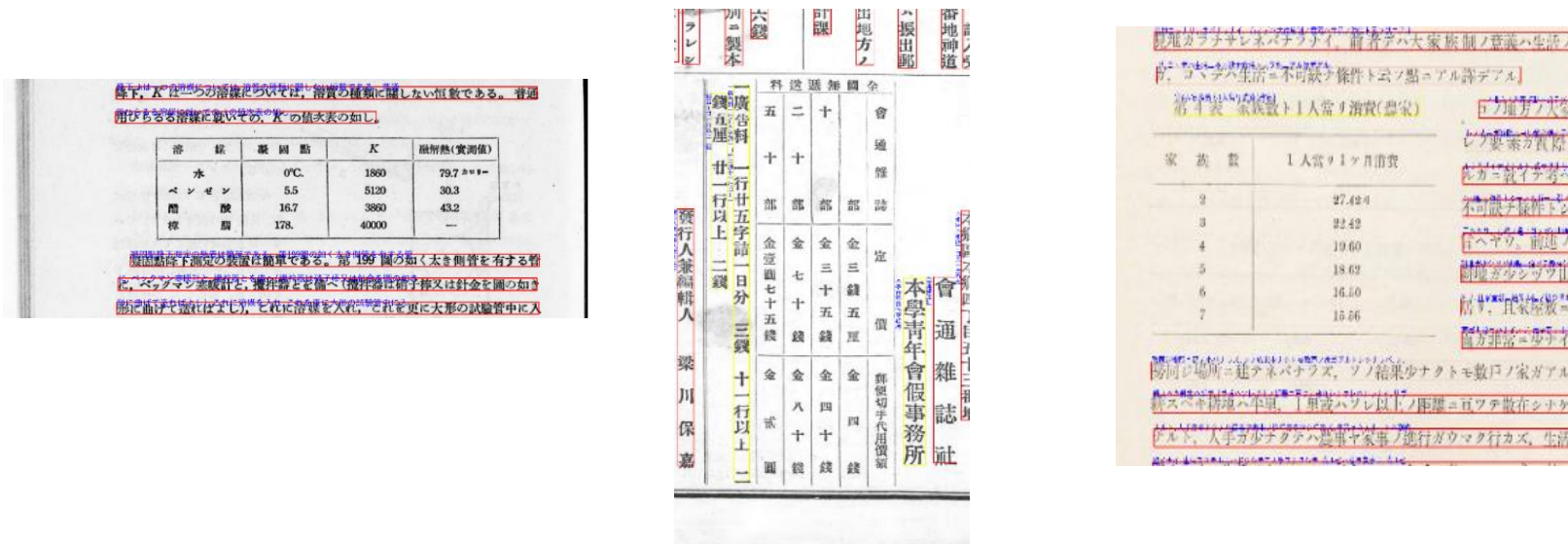
タイトル、ページ数の認識欠落のサンプル（本文ではない柱やノンブルは認識対象外としていた）



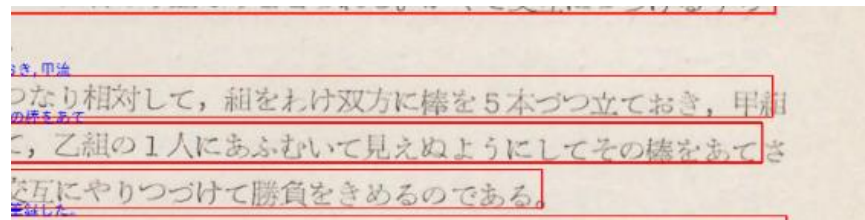
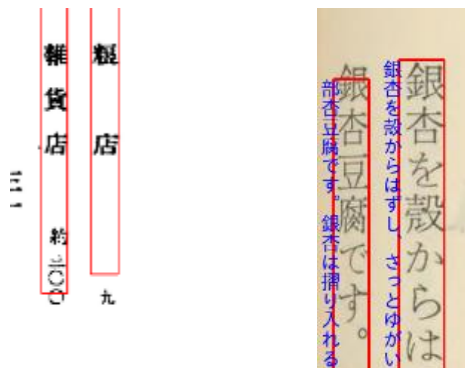
中間報告時点での課題

2. レイアウト、文字間違いの認識性能改善（前ページの続き）

表内文字の認識欠落のサンプル（表内文字は認識対象外としている）

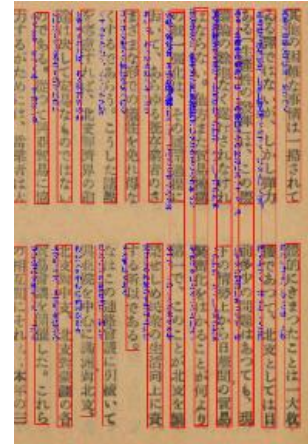


枠の認識欠落（枠が足りていない）のサンプル

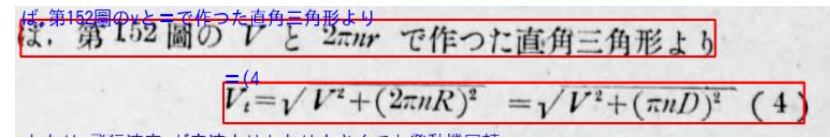
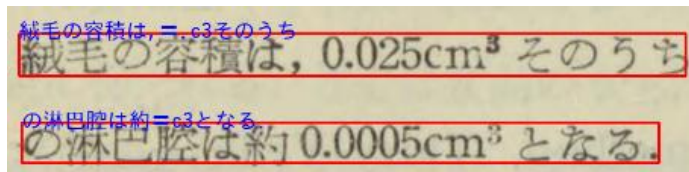
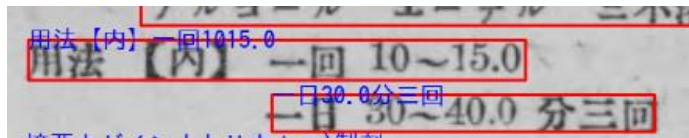


中間報告時点での課題

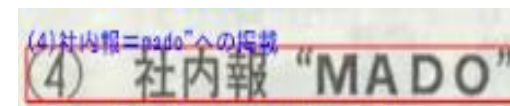
枠の重複認識のサンプル



数字・記号間違いのサンプル



特に150°以上にて溜出する部分のみを採
特に 150° 以上にて



中間報告時点の課題に対する取り組み

1. 性能改善

1) 特に理系文書での欧文、数値、記号などの認識性能改善

→ アノテーションで二になっているインライン欧文に高精度英語OCRでアノテーション

2) レイアウトの縦書き横書き混在の誤認、行の重複の誤認、行の長さの不足

→ ルールベースの後処理で効果を検証、(学習データの増加でも改善)

3) 行認識での二の取り扱い

→ 頻出するアノテーションが二になる文字があればデータセットの修正を検討する想定だったが今回は実施せず

→ 行認識で二が第一候補として認識された場合に第二候補以降を採用するようなオプションを追加

4) 柱・ノンブル対応

→ 柱とノンブルがテキスト化の対象外だったが、オプションで対象/非対象を切り替えられるよう対応

2. 速度改善

1) 個別の手法ごとに現状分割0.6s、傾き補正2.6s、レイアウト認識0.2s、行認識約2.0s、合計5.4s (1ページ辺り2.7s)

速度要件は見開き左右の1ページで2s以内

→ 特に処理時間のかかっている傾き補正の改善を実施

1. 性能改善

1) 理系文書における欧文、数値、記号などの認識性能改善

■ アノテーションで■になっているインライン欧文に対して高性能な英語OCRでアノテーション実施

TrOCR : Transformerを利用した高精度なOCRモデル

精度変化 : 0.818 → 0.825 (欧文が含まれる行に限定した場合 0.730 → 0.972)

欧文アノテーションサンプル

(before)

<LINE DIRECTION="横" TYPE="キャプション" STRING="■" X="527" Y="485" WIDTH="804" HEIGHT="42">

(after)

<LINE DIRECTION="横" TYPE="キャプション" STRING="CORUNDUM AND CARBORUNDUM WHEELS." X="527" Y="485" WIDTH="804" HEIGHT="42">

(before)

<LINE DIRECTION="横" TYPE="キャプション" STRING="■" X="788" Y="550" WIDTH="235" HEIGHT="38">

(after)

<LINE DIRECTION="横" TYPE="キャプション" STRING="FOR LATHE" X="788" Y="550" WIDTH="235" HEIGHT="38">

B 115 店 賣 販 造 製 品 要 科 齒 川 小

CORUNDUM AND CARBORUNDUM WHEELS.
FOR LATHE

改善例

調査=mechle
情を明らかにした調査(Mechler
に関する進展中
ある。室内空気質に関する進展中

した調査(mechler,
苦情を明ら
気質に関する進展中
がある。室

揮発油=或は=油=
揮発油(Volatile oil)或は軽油(Light oil)
て溜出するもの
乃至 300°C にて溜出するもの
燈用石油(Lamp oil)
出するもの...
以上にて溜出するもの
重油(Heavy oil)

揮発油(volatile oil)或は軽油(light oil)
揮発油(Volatile oil)
て溜出するもの.....
乃至 300°C にて溜出
燈用石油(Lamp oil)
出するもの.....
以上にて溜出する
重油(heavy oil)
重油(Heavy oil)

8. 1926
(D. med. W. Jg. 52, Nr. 2, Jan
ライプツィグ

8. 1926
(D. med. W. Jg. 52, Nr. 2, Jan
ライプツィグ

1. 性能改善

2) レイアウトの縦書き横書きの混在、行の重複の誤認、行の長さの不足

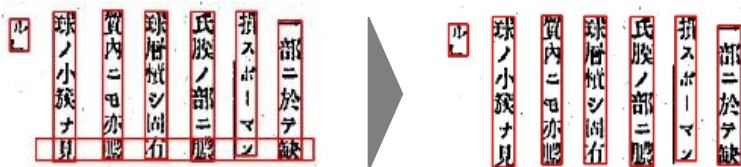
■ ルールベースの後処理を追加して効果を検証

縦横混在：横(縦)行に対して複数の縦(横)行が重なる場合は削除

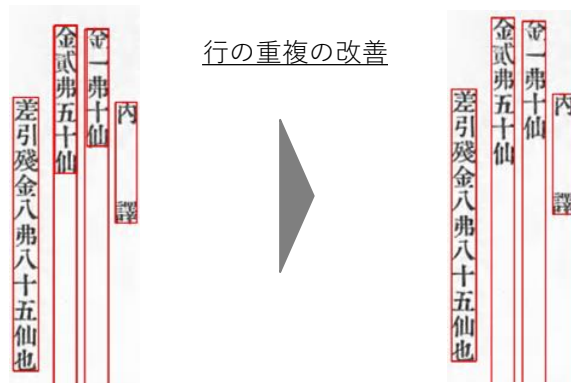
重複：すべての矩形に対して内包されているかをチェックし、内包されている矩形を削除(多少のはみ出しも許容)

長さ不足：学習データの増加で改善

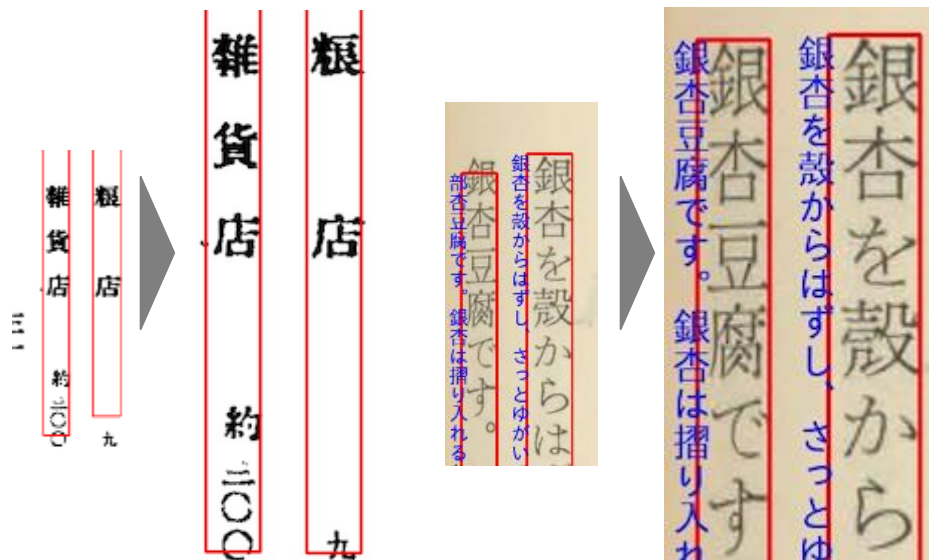
縦横混在の改善



行の重複の改善



長さ不足の改善



2. 速度改善

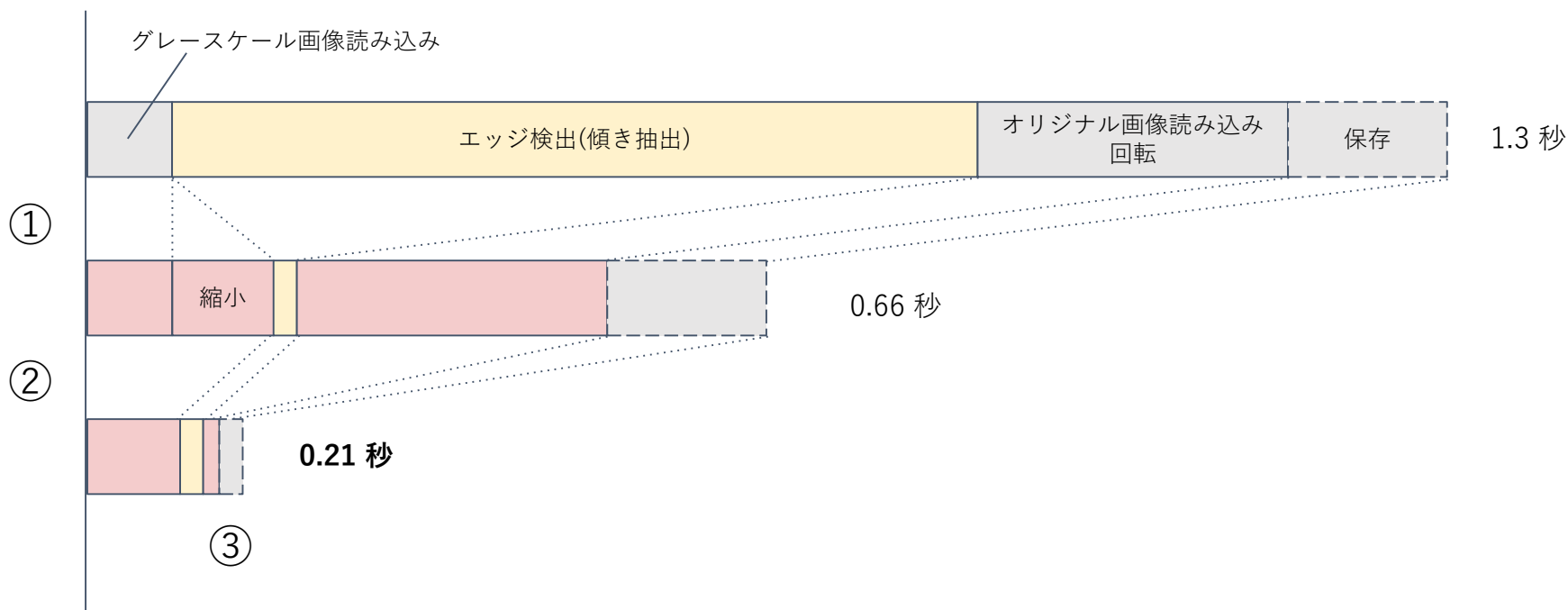
1) 傾き補正と統合プログラムでの速度改善

■ 傾き補正の速度改善

- ① 処理時間の半分以上を占めるエッジ検出(傾き抽出)の時間を入力画像を縮小することで短縮
画像縮小による精度劣化を低減するためにパラメータを調整
- ② 画像の読み込み、縮小や回転などに使用していた画像処理ライブラリをscikit-imageからOpenCVに変更

精度変化：評価画像99枚 x 2(左右) で平均誤差 0.259度 → 0.265度

- ③ 統合プログラムで画像の保存を介さずメモリ上での画像の受け渡しを実行



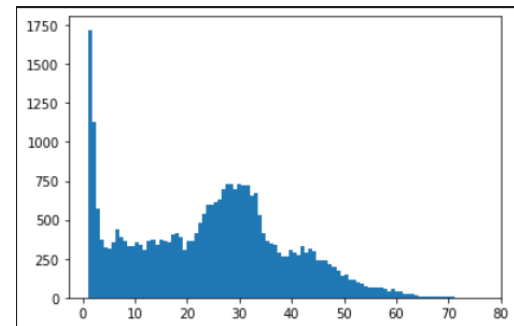
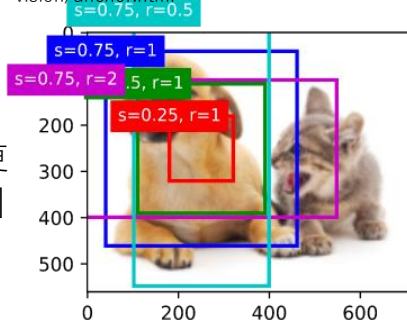
その他の改善検討

■ レイアウト

アンカーボックスサイズ最適化

アスペクト比の構成を全データの矩形情報から以下の通り変更
[1/24, 1/16, 1/8, 1, 8, 16, 24] → [1/32, 1/16, 1/4, 1, 4, 16, 32]

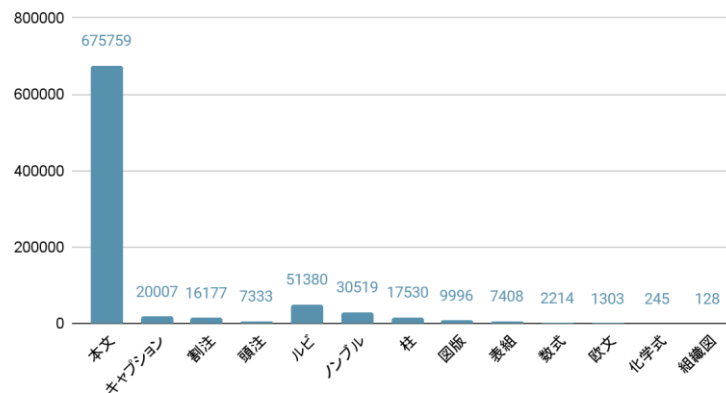
https://d2l.ai/chapter_computer-vision/anchor.html



Equalization Loss

全データの8割を占める本文矩形に偏った学習となるのを回避
少数カテゴリでは一定割合で多数カテゴリデータを無視する

カテゴリごとの矩形の数



■ 行認識

合成漢字データの利用

認識対象文字7083字に対して784字はデータセットに出現しない

Image- True Pred
退津顔端耀暨 c 滔夸他象苦巧犁綫耀董倅幢熙瀨旋杼臺：串鶯燒美海癰歹彡煉睨路包哩墻稼
退津顔端耀暨 c 滔夸他象苦巧犁綫耀董倅幢熙瀨旋杼臺：串鶯燒美海癰歹彡煉睨路包哩墻稼
退津顔端耀暨 c 滔夸他象苦巧犁綫耀董倅幢熙瀨旋杼臺：串鶯燒美海癰歹彡煉睨路包哩墻稼
精度変化：0.788 → 0.800 (出現数50以下の文字でF1 0.267→0.405)

処理結果サンプル

はさるゝ心地めゑ正成然ハ又我も一備して寄手の暇を
 獲させんと、商人形を人々三十ばかり作らせて甲冑を
 せ兵杖を持せ夜中、城のまゝ立せ置前、積板を突並べ其
 後、上勝れたる兵五百人を交て夜、白々と明る、城の中
 り同時、門をどつと開けられ四方の寄手、城内の者死
 者狂ひ、討出たり取圍んで一人も洩さず打取せて我先
 ど攻登り、城の中なる門の聲、目的、矢、誤と作りて射立し
 るぞ城兵、巧みしく暫時の間、數千本の矢を得て人形
 計り本隊は、殘し置兵士の次第、一、城中へ引入し、寄手
 の軍勢、斯とも知らず、城の暗るゝ、お國の商人形を當の軍勢
 ど心得て是を討取んと近付、正成の所存の如く、城を破り
 寄一時、大石を投出しければ、敵兵三百餘人、矢處、打殺さ
 れ半死半生の者五六百人、及びり軍勢、りても、城兵、更、引
 入され、バ寄手、側へ立寄、敵、見るゝ、敵、あて作たる人形、なり
 然すれバ、先刻、是を討んと近付、大石、打殺されたる者、共
 人形、形と討死、たるなり、又是、危、進み、待さうし者、

大勝なりと商人の物、あひひ、成たりける、是より後、
 々々合戦を止め、只徒ら、城を守りて、爲を、棄、なし、寄手の
 大軍、諸病、神は、誘引、れ、城、攻、を、止、め、ける、ゆゑ、城、内、の、者、屈、し
 たるを見て、正成、亦々、一計、を、思、ひ、付、忍、び、の、積、板、持、太、郎、を、招
 け、汝、の、狂、歌、を、好、む、付、斯、の、如、く、狂、歌、を、詠、て、一、つ、の、札、を、
 大、佛、の、陣、前、へ、寄、置、置、て、歸、る、べ、し、と、命、じ、ける、に、板、持、
 是、つ、て、其、如、く、取、計、を、取、り、
 是、所、の、み、見、て、や、い、み、な、り、つ、ら、さ、の
 たり、の、山、の、み、ね、の、く、す、の、木
 と、書、付、て、敵、の、大、佛、の、陣、前、へ、懸、置、け、れ、バ、數、方、の、寄、手、此
 歌、を、見、て、安、ら、ず、思、ひ、疎、さ、ら、大、佛、の、無、念、背、體、お、懸、し、け、れ
 共、爲、へ、は、様、な、く、城、を、白、眼、で、居、た、り、け、り
 ○早、瀬、右、衛、門、忠、義、を、圖、を、案
 丹、波、川、式、部、矢、文、の、事
 夫、國、勢、よ、し、て、都、府、深、き、者、ハ、必、ぞ、人、に、謀、ら、れ、恥、辱、を、取、る、事
 多、し、と、愛、は、足、利、治、部、大、佛、高、氏、ハ、或、時、我、陣、所、ヨ、宗、徒、の、人、々

淡堂云日
 本軍將二
 人突兀入
 陣未遽可
 知矣後脱
 冒始知義
 經一鐵浮圖
 摸得極妙

ニ由ナシト雖、此山師ト以前滿洲チ支配シタル彼ノ軍將
 トハ自カラ別人タルコニ注意セザルベカラ
 此事ハ日記中「サチヤン」城ノ側街ニ二十年以上住居シタル
 支那山丹人「ウチンチン」ノ話チ記載シタル章ヲ見レバ益々
 瞭然タルニ至ルベシ其語ニ曰ク土人ノ口碑ニ從ヘバ昔時
 日本軍將二人アリ此國ニ來レリ其名チ金鳥諸寛永ト云フ
 蘇城ハ其一人ノ建築シ、モノナリ然レハ二人中此地ニ先着
 シタルハ金鳥諸ナルヤ將タ寛永ナルヤ又タ其到着シ、ハ
 何時ナルヤ之ヲ知ルニ由ナシ寛永ハ此地ノ君長トナリ子
 孫世襲スルヲ殆ソト三百年碑アリ銘チ彫刻シ今尙「サチヤ
 シ」城内ニ有シ金鳥諸ノ娘モ亦タ城壁チ築キ「タンキン」城
 ト名附ケタリ金鳥諸及娘ノ權ハ尙「ムタ河」ノ邊ニ保存ス云