

デジタル化画像の活用事例

深層学習手法を用いた 古書籍画像の可読性向上

国立情報学研究所 特任准教授

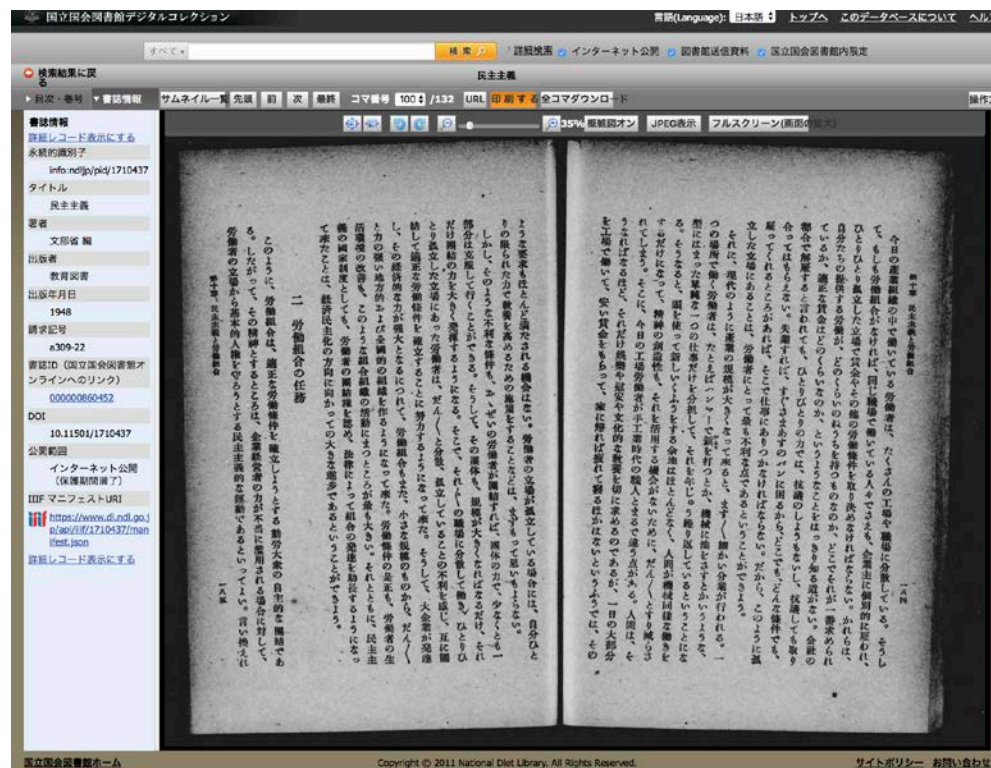
国立国会図書館電子情報部次世代システム開発研究室 委嘱研究員
特定非営利活動法人連想出版 理事

阿辺川武

2018年10月30日

この研究のきっかけ (1/3)

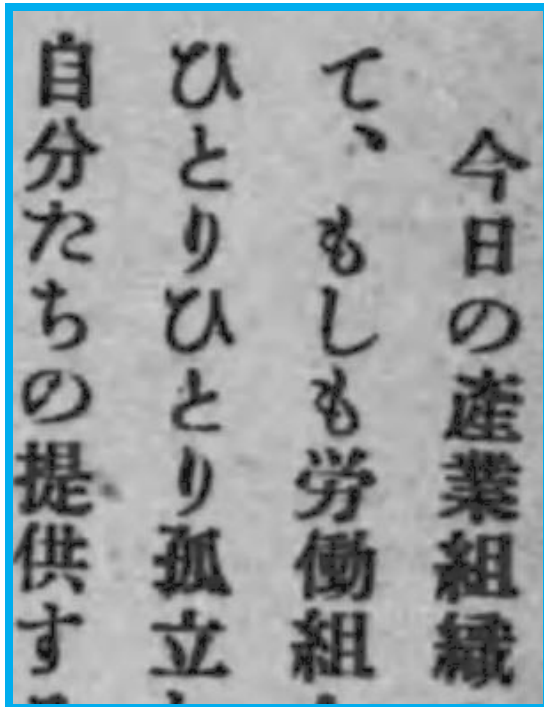
- NDLデジタルコレクション
 - 著作権保護期間が満了したが切れた古書籍が多数公開



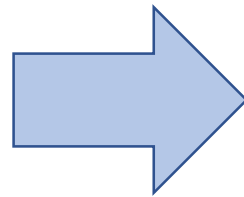
写真撮影の後、
マイクロフィルムで
保存されているため
画像が美しくなく
可読性が悪い

この研究のきっかけ (2/3)

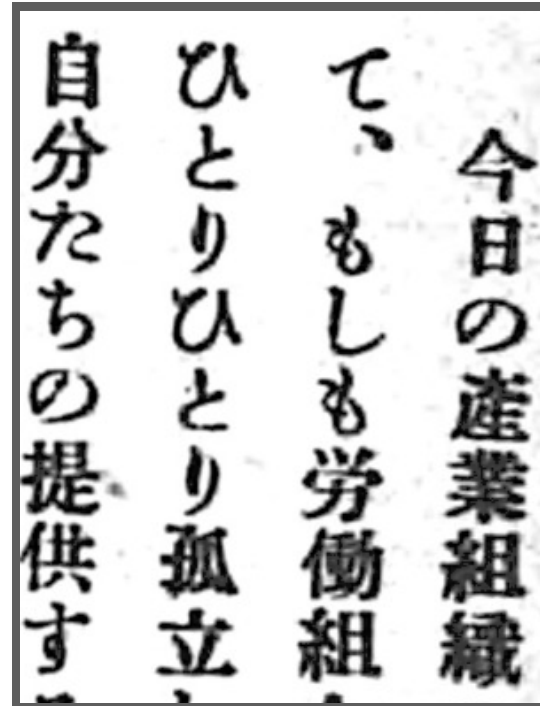
- きれいな画像で書籍を読みたい
 - コントラストや明るさを調整して背景を白色化



元文書画像



各種レベル
補正



補正後画像

背景にゴミが
残ったり、
文字のカスレ
がそのまま

この研究のきっかけ (3/3)

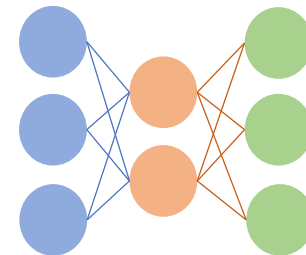
NDL所蔵古書POD (プリントオンデマンド)

The screenshot shows the Amazon.co.jp product page for the book "Philosophy and Science" (哲学と科学との間) by Tanihara Gen (田辺元). The page is in Japanese and features the following elements:

- Header:** Amazon.co.jp logo, navigation links (お届け先, カテゴリー, Amazonポイント, etc.), and a location dropdown set to Tokyo.
- Banner:** A promotional banner for Kindle Store, stating that users can get up to 10% off by buying 5 or more books.
- Product Title:** "Kindleストアでは、哲学と科学との間 (NDL所蔵古書POD [岩波書店]) を、Kindle無料アプリまたはKindle電子書籍リーダーで今すぐ読むことができます。" (On the Kindle Store, you can read "Philosophy and Science" (NDL Collection Old Book POD [Rinsen Shoten]) with the Kindle free app or Kindle e-reader right now.)
- Book Cover:** A small image of the book cover, which is a simple, light-colored design with the title and author's name.
- Pricing:** The Kindle version is priced at ¥324. The print-on-demand (POD) version (ペーパーバック) is priced at ¥2,673. A note indicates that the POD version is a reprint of a 2015/2/27 edition.
- Special Offer:** A banner for a special campaign, stating that users can get 30% off on daily necessities by buying 1 or more items.
- Delivery Information:** A note stating that the book can be delivered by 4/17 (Sunday) if ordered by 8:00 PM on 4/16.
- Points:** A banner indicating that users can earn 300 points by purchasing the book.

<https://www.amazon.co.jp/dp/4802005970/>

CycleGAN



- ディープラーニングを用いた画像変換手法の 1 つ (Yan, 2017) <https://github.com/junyanz/CycleGAN>
 - 馬をシマウマに
 - 夏景色を冬景色に
 - 写真をモネ風、セザンヌ風、ルノワール風に
 - オレンジをりんごに



horse → zebra

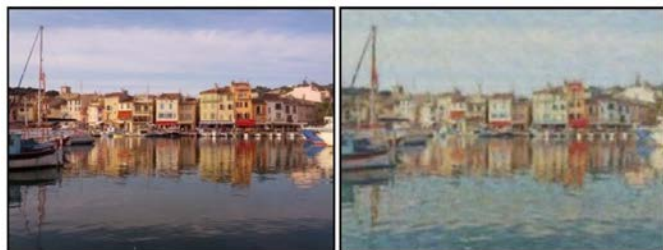


photo → Monet



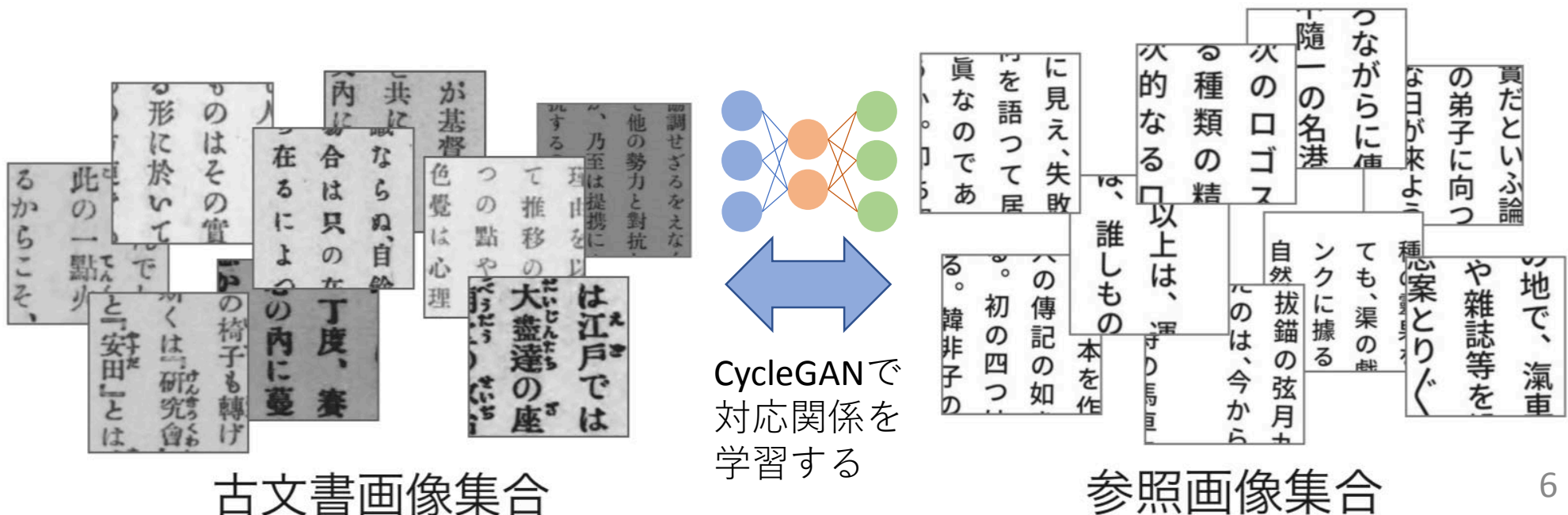
summer → winter



orange → apple

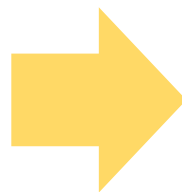
フォントスタイル変換

- CycleGANを古書籍の文書画像に応用できないか？
- 2つの画像集合を用意し、対応関係を機械学習させる
 - 古書籍の文書画像集合：画像から切り出す
 - きれいな文字の画像集合(参照画像集合)：組版ソフトで作成



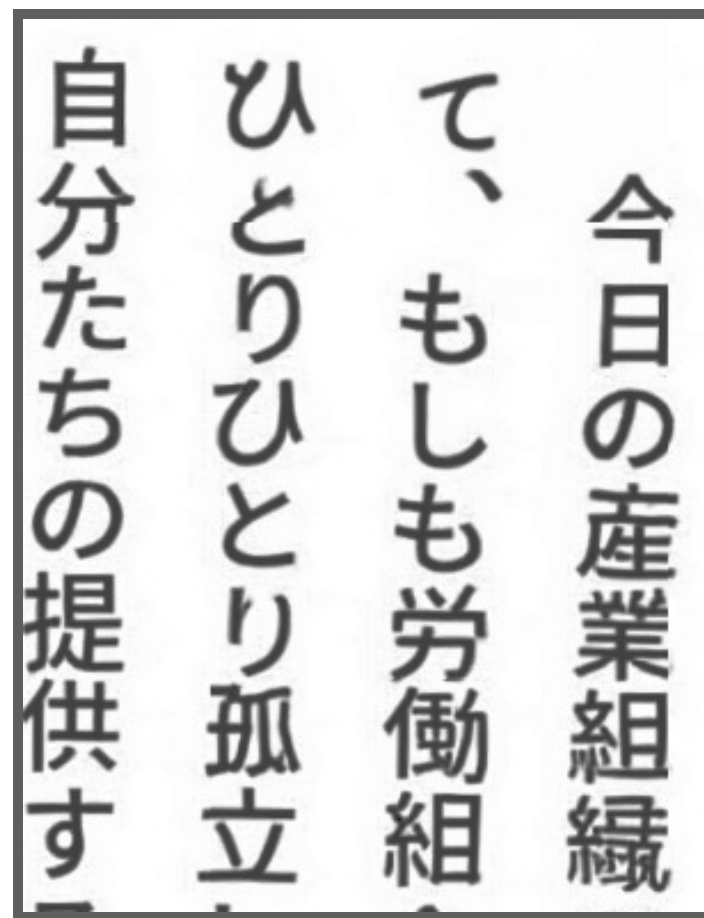
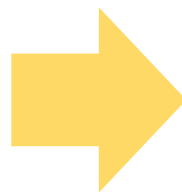
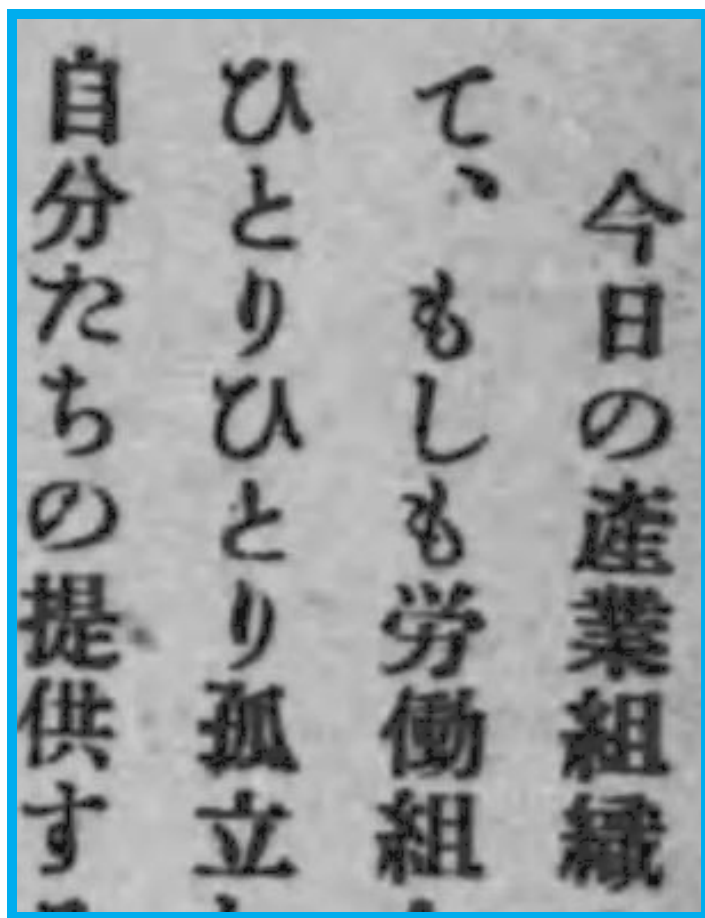
結果（参照画像に明朝体）

今日の産業組織
て、もしも労働組
ひとりひとり孤立
自分たちの提供す

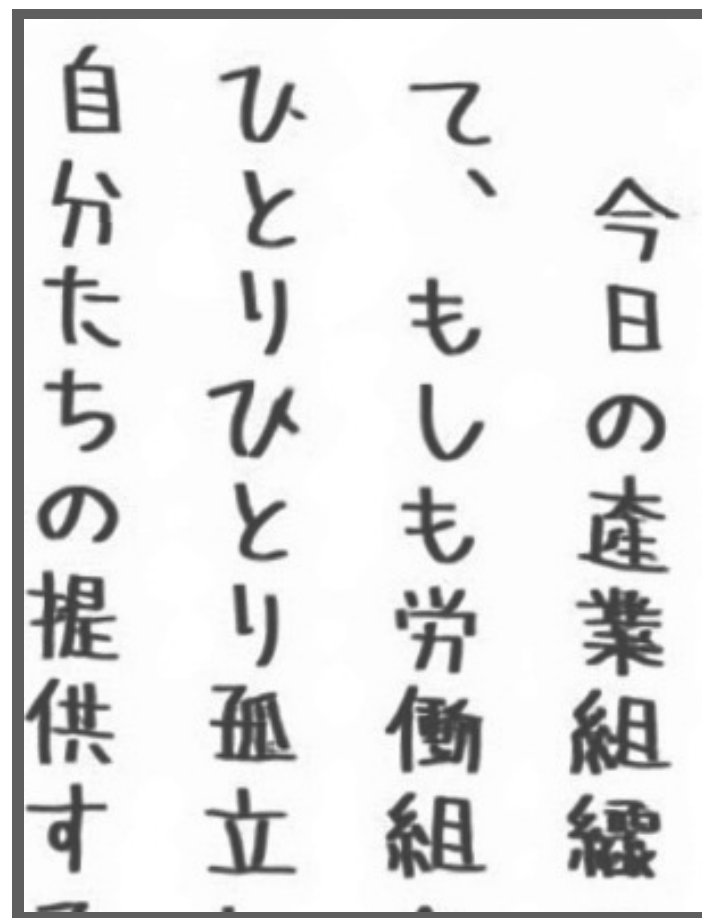
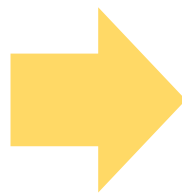
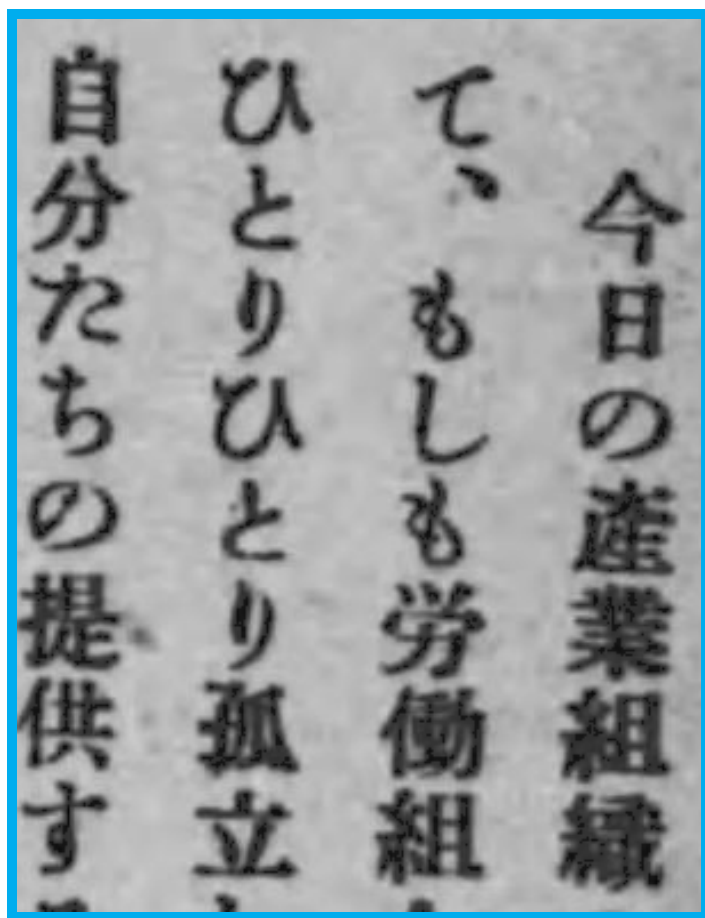


今日の産業組織
て、もしも労働組
ひとりひとり孤立
自分たちの提供す

結果（参照画像にゴシック体）



結果（参照画像に手書きフォント）



今日の産業組織の中で働いている労働者は、たくさん工場や職場に分散している。そうして、もしも労働組合がなければ、同じ職場で働いている人々でさえも、企業主に個別的に雇われ、ひとりひとり孤立した立場で賃金やその他の労働条件を取り決めなければならない。かれらは、自分たちの提供する労働が、どのくらいのぬうちを持つものなのか、どこでそれが一番求められているか、適正な賃金はどのくらいなのか、というようなことをはっきり知る道がない。会社の都合で解雇すると言われても、ひとりびとりの力では、抗議のしようもないし、抗議しても取り合ってもらえない。失業すれば、すぐさまあすのパンに困るから、どこでも、どんな条件でも、雇ってくれるところがあれば、そこで仕事にありつかなければならない。だから、このように孤立した立場にあることは、労働者にとって最も不利な点であるということができよう。

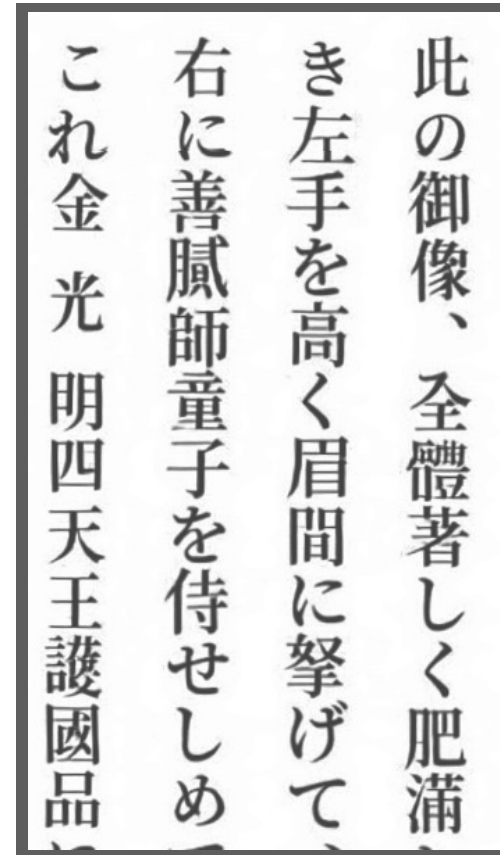
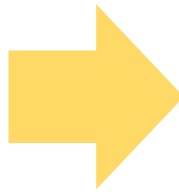
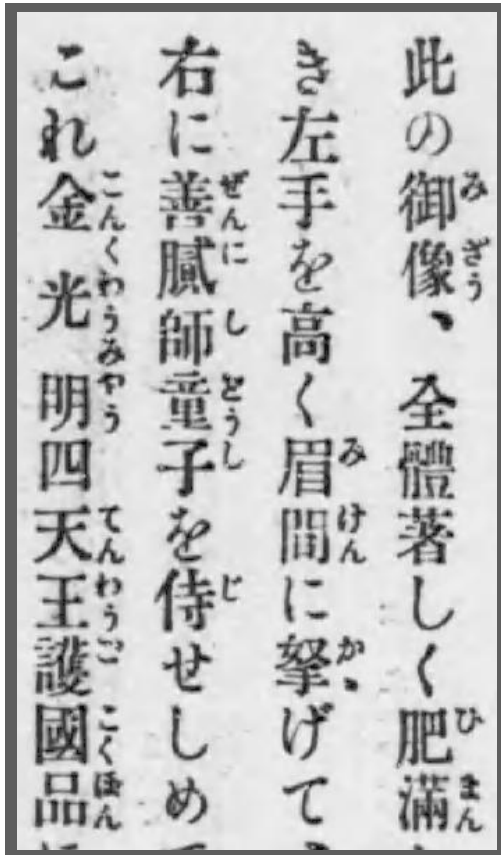
それに、現代のように産業の規模が大きくなって来ると、ますます細かい分業が行われる。つの場所で働く労働者は、たとえばハンマーで鋸を打つとか、機械に油をさすとかいうような型にはまった単純な一つの仕事だけを分担して、それを年じゅう繰り返しているということになる。そうになると、頭を使って新しくいふうをする余地はほとんどなく、人間が機械同様な働きをするだけになって、精神の創造性も、それを活用する機会がないために、だん／＼とすり減らされてしまう。そこに、今日の工場労働者が手工業時代の職人とまるで違う点がある。人間は、そうしなければならぬほど、それだけ娯楽や慰安や文化的な教養を切に求めるのであるが、一日の大部分を工場で働いて、安い賃金をもらって、家に帰れば疲れて寝るほかはないというふうでは、その

今日の産業組織の中で働いている労働者は、たくさん工場や職場に分散している。そうして、もしも労働組合がなければ、同じ職場で働いている人々でさえも、企業主に個別的に雇われ、ひとりひとり孤立した立場で賃金やその他の労働条件を取り決めなければならない。かれらは、自分たちの提供する労働が、どのくらいのぬうちを持つものなのか、どこでそれが一番求められているか、適正な賃金はどのくらいなのか、というようなことをはっきり知る道がない。会社の都合で解雇すると言われても、ひとりびとりの力では、抗議のしようもないし、抗議しても取り合ってもらえない。失業すれば、すぐさまあすのパンに困るから、どこでも、どんな条件でも、雇ってくれるところがあれば、そこで仕事にありつかなければならない。だから、このように孤立した立場にあることは、労働者にとって最も不利な点であるということができよう。

それに、現代のように産業の規模が大きくなって来ると、ますます細かい分業が行われる。一つの場所で働く労働者は、たとえばハンマーで鋸を打つとか、機械に油をさすとかいうような型にはまった単純な一つの仕事だけを分担して、それを年じゅう繰り返しているということになる。そうになると、頭を使って新しくいふうをする余地はほとんどなく、人間が機械同様な働きをするだけになって、精神の創造性も、それを活用する機会がないために、だん／＼とすり減らされてしまう。そこに、今日の工場労働者が手工業時代の職人とまるで違う点がある。人間は、そうしなければならぬほど、それだけ娯楽や慰安や文化的な教養を切に求めるのであるが、一日の大部分を工場で働いて、安い賃金をもらって、家に帰れば疲れて寝るほかはないというふうでは、その

ルビの削除

- ルビが全く存在しない参照画像集合を用いるとルビが消える



OCRの精度
向上に
寄与する

親と頼んで隠れつゝ、一心讀誦、毘沙門經、咽も裂けよと叫ぶ折から、ア、ラ不思議や、大木自らにどうと仆れて、三つの妖鬼は其の下敷、板の如くに潰されて、これは又見るも無慚、敢へなき最期を遂げてしまつた——こゝから最急行。其の翌朝伊勢人登山して、峯延と師檀を約す、鞍馬山寺の隆運、こゝに源を發するといふ。

閑話休題、鞍馬安置の毘沙門天は、伊勢人の手に勸請した靈像。——尤も藤原前期との説もある——都南門樓の兜跋毘沙門に對して、鞍馬は正しく帝京の坤位、王城鎮守の靈山だ。こゝに北方天の靈體を安置して、千代萬代の末永く、皇基を守護し、都城を衛護せしめむとは、忠臣伊勢人の、夢寐なほ忘れ得ぬ宿題であつた。さればこそ此の御像、全體著しく肥滿して強剛に、分けても顔面の緊張物々しく、寶塔を捧ぐべき左手を高く眉間に擎げて、四六時中皇都俯瞰の相丰生けるが如く、且又左に吉祥天、右に善膩師童子を侍せしめて、天王の威光更に一段を加ふ、而して兩脇侍立の形式は、これ金光明四天王護國品に説けるところ、伊勢人の意を用ふる、まことに厚く深きものがある。

親と頼んで隠れつゝ、一心讀誦、毘沙門經、咽も裂けよと叫ぶ折から、ア、ラ不思議や、大木自らにどうと仆れて、三つの妖鬼は其の下敷、板の如くに潰されて、これは又見るも無慚、敢へなき最期を遂げてしまつた——こゝから最急行。其の翌朝伊勢人登山して、峯延と師檀を約す、鞍馬山寺の隆運、こゝに源を發するといふ。

閑話休題、鞍馬安置の毘沙門天は、伊勢人の手に勸請した靈像。——尤も藤原前期との説もある——都南門樓の兜跋毘沙門に對して、鞍馬は正しく帝京の坤位、王城鎮守の靈山だ。こゝに北方天の靈體を安置して、千代萬代の末永く、皇基を守護し、都城を衛護せしめむとは、忠臣伊勢人の、夢寐なほ忘れ得ぬ宿題であつた。さればこそ此の御像、全體著しく肥滿して強剛に、分けても顔面の緊張物々しく、寶塔を捧ぐべき左手を高く眉間に擎げて、四六時中皇都俯瞰の相丰生けるが如く、且又左に吉祥天、右に善膩師童子を侍せしめて、天王の威光更に一段を加ふ、而して兩脇侍立の形式は、これ金光明四天王護國品に説けるところ、伊勢人の意を用ふる、まことに厚く深きものがある。

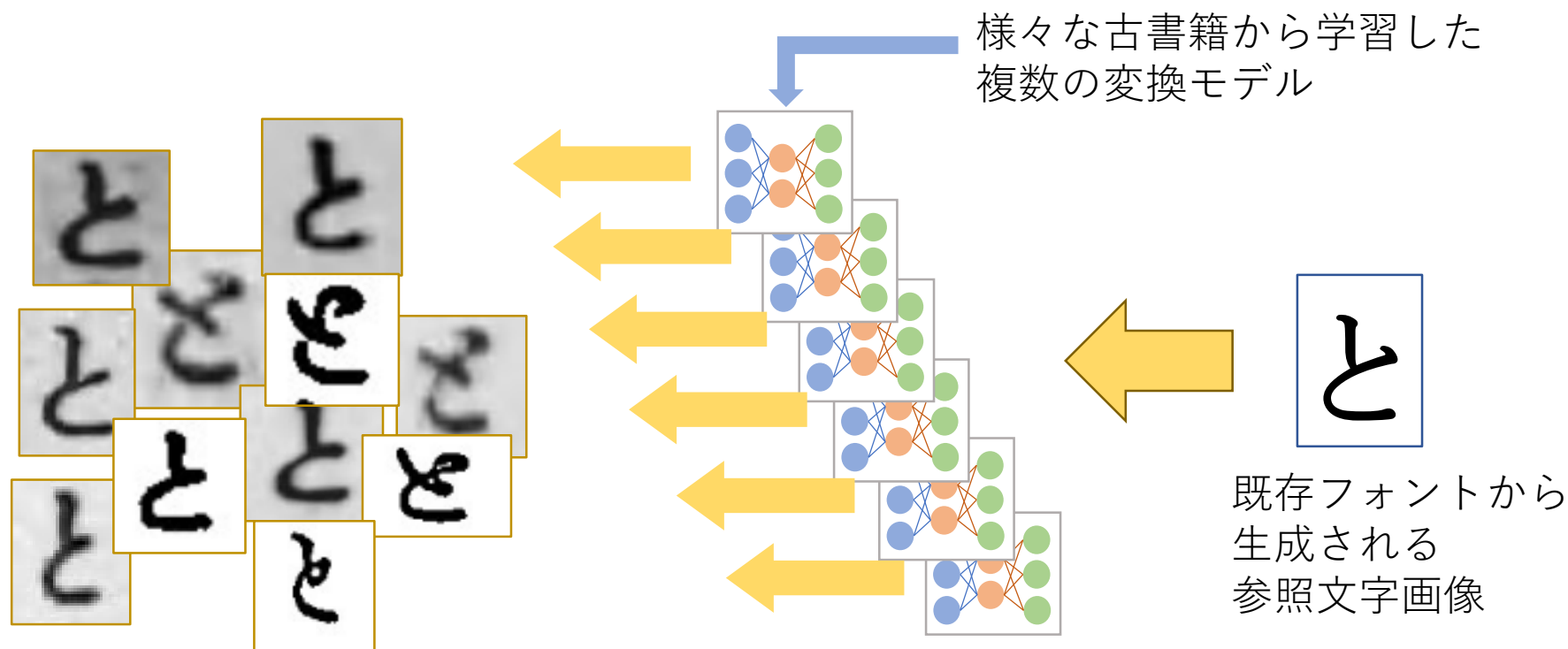
OCRの精度

フォントスタイル変換前と変換後でOCR精度を比較

	1710437-100		983386-156 (ルビあり)	
	無修正	変換後	無修正	変換後
OCR A	0.949	0.969	0.465	0.922
OCR B	0.898	0.961	0.860	0.938
OCR C	0.982	0.981	0.903	0.925

今後の予定

- 古書籍に特化したOCRの作成
 - フォントスタイル変換手法を使って、
全ての文字の古文書画像を自動生成する



参照画像から古書籍画像へ逆変換（いままでとは逆）