#### 第8回日本語の歴史的典籍国際研究集会

# 「みをつくし」プロジェクト: AIくずし字認識研究の展開

#### カラーヌワット・タリン

Senior Research Scientist Google Research, Brain team

発表内容は個人の見解であり、所属組織を代表するものではありません。







#### 『宇津保物語』(国文学研究資料)



日本語の歴史的典籍の国際共同研究ネットワーク構築計画

(略称:歷史的典籍NW事業)

歴史的典籍22万6千点、コマ数では 2400万コマを撮影した。今後30万点の 画像を公開する予定。



#### くずし字データセット

江戸時代の古典籍 44 冊

6151 ページ

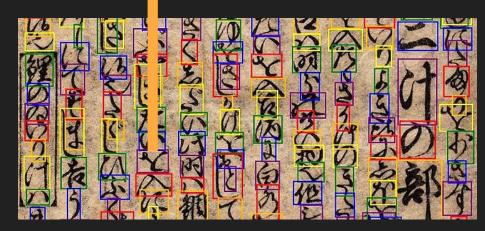
1,086,326 文字

4,328 文字種

Unicode	Image	Х	Υ	Block ID	Char ID	Width	Height
U+81EA	100241706_00004_2	1852	1736	B0001	C0001	104	219
U+5E8F	100241706_00004_2	1816	2096	B0001	C0002	152	296
U+82E5	100241706_00004_2	1465	951	B0001	C0003	172	216
U+3044	100241706_00004_2	1495	1218	B0001	C0004	143	69
U+6642	100241706_00004_2	1465	1338	B0001	C0005	168	179
U+306E	100241706_00004_2	1497	1567	B0001	C0006	123	152
U+6C17	100241706_00004_2	1504	1754	B0001	C0007	145	200
U+5F37	100241706_00004_2	1479	2034	B0001	C0008	163	145
U+306B	100241706_00004_2	1521	2239	B0001	C0009	83	237

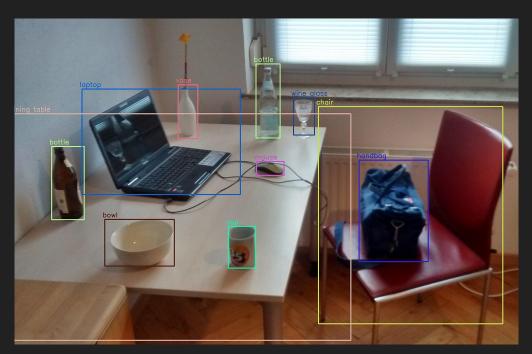
Unicode, x, y, w, h

を U+3029, 512, 418, 56, 47

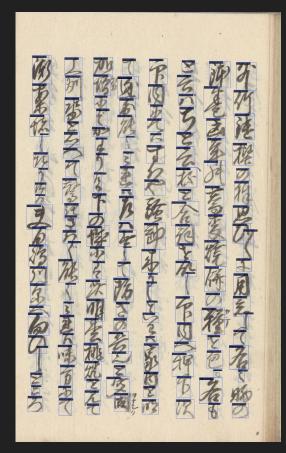


国文学研究資料館&国立国語研究所作成

### くずし字認識の手法:物体検出の手法



Wikipedia



### くずし字認識の手法:物体検出の手法

やける小石川御門内するかだい小川 やけるす きやばより 呉服ば内大名 御馬やかし迄無事す せんぼり七曲泉ばし新橋浅草御門 松平肥後守様上中やき共やけ かぢ丁 転 十一中橋南てんま丁二丁目横丁 本心舟後様其 中やしき同上やしきやけるときはじら 小路 前 後御大名少々破 多 腰掛やける同南 北こんや丁白 急やしき 同東 すだ町ゟ 今川ばし日本 ば 上 家 柳ずし少々破 損 同所北の方御蔵前 十一橋御門外 ごぢいん原 きじば 中ばまで家蔵 筋 違御毘少々破損外 神田さみ 八日河岸上やける土蔵 此間に 辰 之口向 森 川出羽守 様 酒井うた様 令 ちこいなば 町家共大破損っき のこる京 しき 柳 損 多く火 火 者 垂 大日本橋 ゟ 和田倉御門内大 し御門内少破損 丁具 野新道五郎三丁過夕日 ばる新 事 外にやしき五六 足丁 番 大破損 橋上御屋 じ小田原丁南 松平下従守様 炭丁村木丁 しかい御門内 所同 田 松平ぞ 安 所 方其 軒

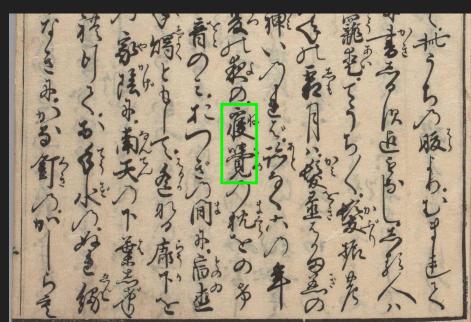
画像:みんなで翻刻プロジェクト

日比谷御門外 少

損

# くずし字認識の手法:物体検出の手法





#### (3) |字認識研究の展開:古活字の研究

Α

きあず人ろもとる

色りのぬあ下刀の成るに

200

) おをもくな

つはう那

そろとなわかりけっなり

かかと

ろんさん

しられる

おからふ

るけっちてを見り

В

Po

わりつのと

かってる

8

きのもけかのいるか

くばきえ病の and a

きあれ人ろもと 我 色力のぬ ノおをもくいなつはりか あべりり 78

de

うなとなわからけっない ろんまていられるのうかからふ いるは

当二

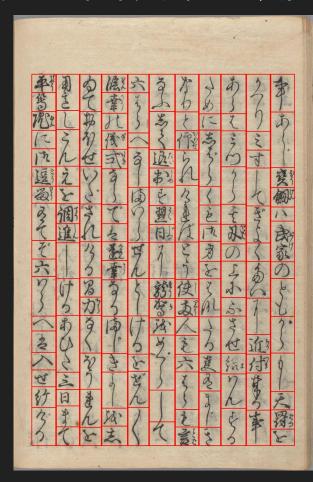
をかと

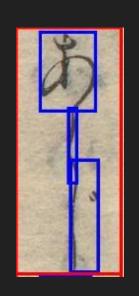
と低きえ家のやとわなりと いくそろ きのもけむのいるも いわるるか n

A + B

ちるこ をかとい 我 色りのぬあて刀りのもは くはきえ前の よ人ろもとる ノおをもくいなつはり わかわけっすり さらけるのうかか 色月 きのもけれるいるか いわるる りつ なるとろ でするる

### くずし字認識研究の展開: 古活字の研究

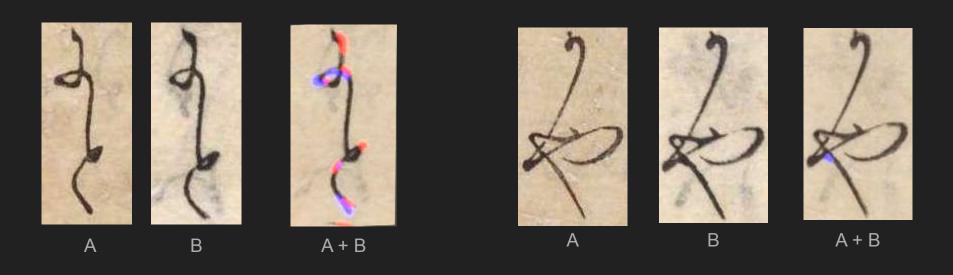








### くずし字認識研究の展開: 古活字の研究

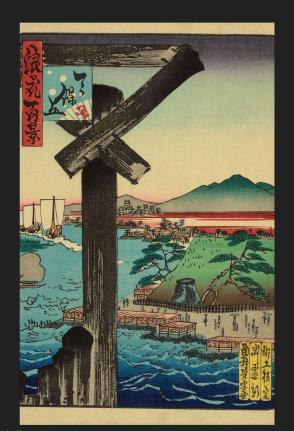


A: 内閣文庫蔵:嵯峨本第一種本『徒然草』 B: 内閣文庫蔵:嵯峨本第四種本『徒然草』

# 「みをつくし」プロジェクト

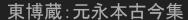
### みをつくし(澪標)

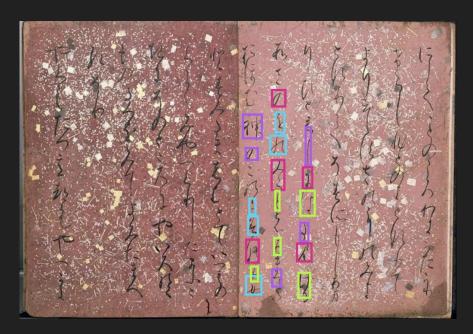
- ■『源氏物語』第14巻「みをつくし」(澪標)。
- 「みを(船の水路)を示すために立ててある杭」の
  意。
- 「身を尽くし」の掛詞。
- 「みをつくし」が人々の水先案内となるように、「みをつくし」プロジェクトがくずし字資料を読むための道案内となることを目指している。

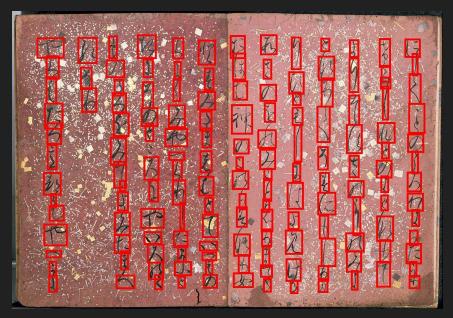


14

# 新AIくずし字認識モデル: RURI







旧モデル RURIモデル



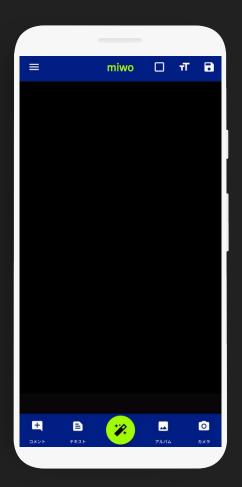


#### 「みを」くずし字認識アプリ

- 画像からくずし字を認識し、現代日本語文字に変換する。
- 2021 年 8 月 30 日にリリース。 現在、10万回以上ダウンロード。
- iOS、Android両方リリースした。Webアプリとデスクトップアプリもリリース可能。
- 現在まで認識した画像の枚数は 100万枚以上。

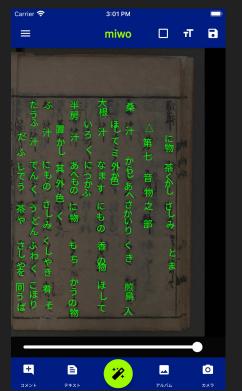


2022年10月7日にグッドデザイン賞(システム・サービス部門)を受賞した。



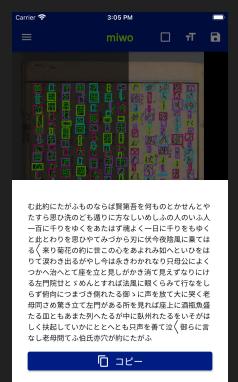


# 「みを」アプリの主な機能

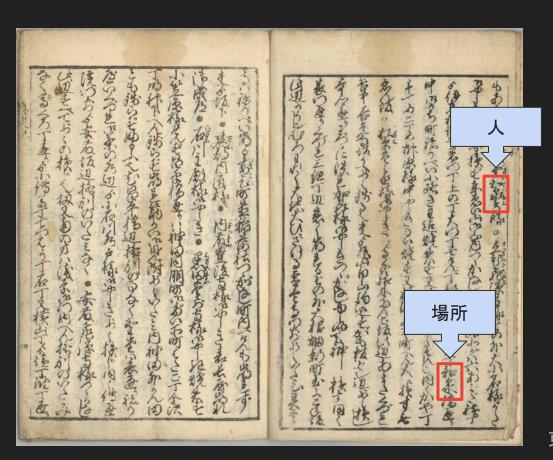








### 「つくし」プロジェクト



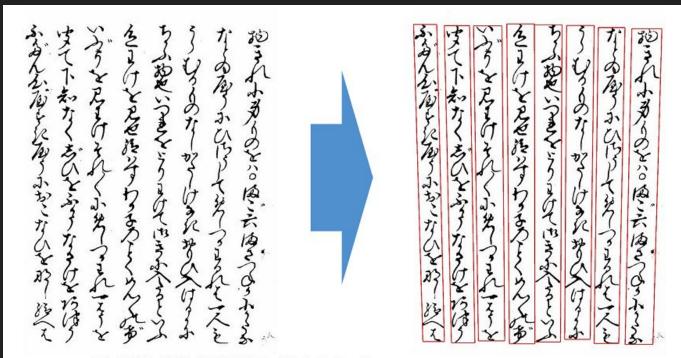
#### 和泉



- ◆ くずし字認識を使って新しい研 究方法を開拓する。
- 検索機能
  - 全文検索
  - ユニバーサル検索
    - 画像
    - 地図
    - ナレッジグラフ
- 事文字情報を分析する方法。
- 古典籍を探し尽くす。

#### くずし字認識の手法:くずし字行領域認識

SIGNATE 凸版印刷株式会社 くずし字認識チャレンジ① くずし字 行領域認識アルゴリズム作成



出典:人文学オープンデータ共同利用センター(http://codh.rois.ac.jp/) 『日本古典籍くずし字データセット』に含まれる『吉利支丹物語』(国文学研究資料館所蔵、doi: 10.20730/200006665、 CC BY-SA(https://creativecommons.org/licenses/by-sa/4.0/)にて配布、新日本古典籍総合データベースより)を加工した画像加工内容: 行領域を4点ポリゴンで囲って赤枠で表示