

## リーディングライフ・ログ

木村 崇志 † Rong Huang †, フォン ヤオカイ †, 内田 誠一 †

岩村 雅一 ‡, 大町 真一郎 \*, 黄瀬 浩一 †

†九州大学 ‡大阪府立大学 \* 東北大学

E-mail: uchida@ait.kyushu-u.ac.jp

### Abstract

リーディングライフ・ログ (Reading-Life Log) とは、我々が日々読んでいる文字を自動的に認識し、記録するシステムである。その実現のためには、読んでいる範囲の画像を取得し、その領域に対して文字認識を行えばよい。このように考え方は比較的シンプルである反面、その実現については幾つかの課題がある。例えば、動画像として連続的に入力される文書画像について、フレーム間のオーバーラップを考慮しながら適宜ログ化する必要がある。さらには文字を認識するために、できる限り高精細かつ動きボケの少ないカメラを頭部に装着する必要がある。本研究では、こうした課題について一定の解決を図った上で、全体的なシステムの実装を行った。また書籍やスマートフォン画面を対象としたログ化実験を行い、その性能を評価した。

### 1 概要

センサ技術の発達により、様々なライフログを気軽に収集できるようになった。例えば、脈拍を記録することで健康状態をログ化でき、また GPS 情報を記録することで、行動経路をログ化できる。さらにこうしたログを解析することで、長時間データに依拠したユーザー行動のモデル化や正常・異常解析等の可能性も広がる。

我々は様々なライフログの対象から、人の読んだ文字の収集に注目した。文字は、我々が生活する上で必要不可欠なものであり、多くの情報を文字から得ている。つまり、個人の読んだ文字を得ることは、非常に有益なことであり、多くの利用法も考えられる。本研究は、以上のことから日常的に読んでいる文字情報の自動獲得を目指し、これをリーディングライフ・ログ (Reading-Life Log) と呼ぶ。

リーディングライフ・ログの利用範囲は多様である。例えば、特定文字がある情景を見た経験があるとし、それをしばらく後に再び思い出したいとする。このとき、その情景画像中の特定文字を認識しておけば、後にそれらの単語を用いて検索するだけで、この情景画像を

呼び出せる。すなわち、画像をキーとした検索とは異なり、文字情報を用いた曖昧性が少ない検索が実現できる。

リーディングライフ・ログは、基本的に、(i) 読んでいる範囲の画像を取得し、(ii) その領域に対して文字認識を行う、という2ステップで実現される。このように基本的考え方は比較的シンプルである。しかしその反面、実現までには幾つかの解くべき課題がある。例えば、(i) については、文字を認識するためにできる限り高精細かつ動きボケの少ないカメラを頭部に装着する必要がある。また目で読んでいる文字の位置も何らかの精度基準の下で、特定する必要がある。(ii) については、動画像として連続的に入力される文書画像について、フレーム間のオーバーラップを考慮しながら適宜ログ化する必要がある。すなわち、オーバーラップ部を2重にログ化する必要はなく、またオーバーラップ部が完全一致していないような場合はより適切な方の認識結果を選択する仕組みも必要となる。

ところで我々は以前、上記 (i) の課題について、アイトラッカで視線情報を抽出し、視点位置の文字を逐次獲得するシステムを作成した [1]。視点位置が確実に読んだ文字に合致していれば、これにより、読んだ文字が一つずつ認識・記録されることになる。しかし、その後更なる高精度化を図るべく実験を重ねたところ、現状のアイトラッカでは、装着者の体の動きやキャリブレーションの誤差により、十分な精度の視線情報を得ることができないことがわかった。さらにアイトラッカに付随するシーン撮影カメラは、解像度が低く、またシャッタースピードもコントロールできないので動きボケも激しいため、その後の文字認識処理には不向きである。

以上の理由により、本稿では、アイトラッカを用いないシステムを作成する。具体的には、First-Person Vision のように、小型の高解像度・高シャッタースピードを持つカメラを、頭部に、かつ顔の向きに一致させて装着した。そしてそのカメラの画像の画面中央付近の画像を取得することとした。完全な視点位置ではないため、厳密には読んでいない文字も含まれる可能性があるが、逆に言えば、読んだ文字を見落とすことは

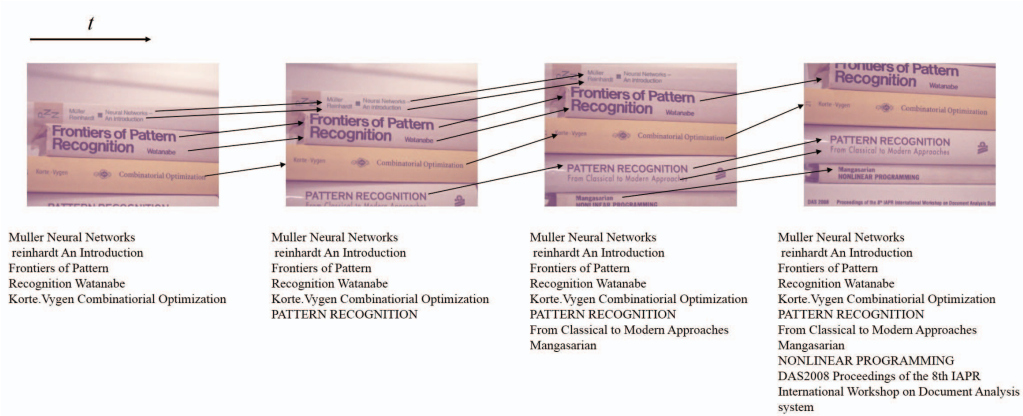


図 1 認識結果統合の例

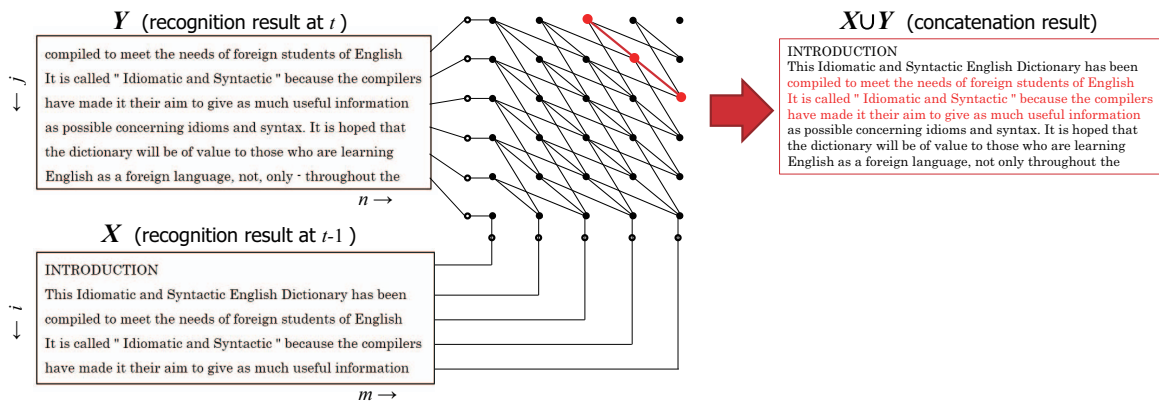


図 2 フレーム間認識結果統合プロセス

まずないと期待できる。

また上記 (ii) の課題については、一種の 2 次元 DP マッチングを用いて、フレーム間の認識結果の 2 次元の整合を図りながら、フレーム間統合を行う。この手法のポイントは、フレーム画像をそのままビデオモザイク処理するのではなく、テキスト（文字認識結果）のモザイク処理となっていることである。すなわち画像レベルでのフレーム間位置合わせ (registration) やトラッキングを行わず、そこで認識されたテキストのみを用いて、フレーム間統合を行う。

本稿ではまず、以上の手法について、その詳細を述べる。次にこのシステムで、屋外、屋内にある看板や文章を対象とした実験を行い、統合結果の精度の確認を行った結果について述べる。

## 2 システムの概要

### 2.1 視界中央の画像の取得

頭部方向の画像の取得には、POINT GREY 社の Flea3 (USB3.0) を利用する。同カメラは小型で頭部装着が可能なサイズながら、 $2080 \times 1552$  画素と高解像度であり、またシャッタースピードも  $1/2000[\text{sec}]$  程度ま

で速くできる（添付ソフトウェア利用）。この速さは、装着者の動きによる画像のモーションブラーを防ぐためには重要である。一方、高速シャッタースピード下では、当然撮影画像は暗くなる。ただしこの問題については、ダイナミックレンジ向上のための輝度変換（実際には 2 値化）をすることで解決でき、実際、その後の認識処理にはほとんど影響ないことを確認している。

### 2.2 認識処理

本稿では、取得した画像に対して、可変閾値処理、大津の 2 値化、平滑化を組み合わせて 2 値化を行い、市販の OCR エンジン (ABBYY 社, FineReader Engine 10) に入力して文字認識を行なっている。同エンジンは、情景内文字検出・認識に特化されたものではないが、それでもある程度のレイアウト解析や文字行抽出、そして情景内文字認識は可能である。例えば、画像中に  $N$  行の文書が含まれていれば、理想的には、その  $N$  行分の認識結果が  $N$  行分の文字列として求まることになる。

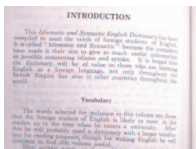
よく知られているように、カメラベースの文字認識は非常に困難な課題である [2, 3]。しかし現状の実装において、我々はまさに上記の 2 値化を施す程度の処理に留まっている。その結果、後の実験に示されるように、誤認識が発生する。ただし、こうした誤認識は、次

INTRODUCTION  
This Idiomatic and Syntactic English Dictionary has been compiled to meet the needs of foreign students of English. It is called "Idiomatic and Syntactic" because the compilers have made it their aim to give as much useful information as possible concerning idioms and syntax. It is hoped that the dictionary will be of value to those who are learning English as a foreign language, not only throughout the British Empire but also in other countries throughout the world.

Vocabulary  
The words selected for inclusion in this volume are those that the foreign student of English is likely to meet in his studies up to the time when he enters a university. After this he will probably need a dictionary with a larger vocabulary for reading purposes, though for writing English he will continue to find this volume useful.

Most archaic words, or those which are likely to occur only in purely scientific and technical contexts, have been excluded. Colloquial and slang words and expressions have been included if they are of the sort likely to be found in books (e. g. modern fiction and drama) read by students. Foreign words and Latin words and phrases of common occurrence in English have also been included.

Definitions  
Definitions have been made as simple as possible. Where definitions in easy, common words was not practicable or satisfactory, pictures and diagrams have been supplied. A lobster, in the Concise Oxford Dictionary is defined as "a large marine stalk-eyed ten-footed long-tailed edible crustacean with large claws formed by the first pair of feet, bluish-black before and scarlet after boiling; it's flesh as food". The foreign student of English, if he is a beginner, is likely to be puzzled by certain words in this definition ("stalk-eyed" and "crustacean", for example). The C. O. D. was not written specially for him. The ordinary user of the C.O.D. is likely to be a person who knows quite well what a lobster is and who refers to this word in a dictionary only when he needs exact and scientific information of the kind given in the above admirably concise and complete definition. The foreign student usually needs only to identify the new



(a) 取得画像例

INTRODUCTION  
This Idiomatic and Syntactic English Dictionary has been compiled to meet the needs of foreign students of English. It is called "Idiomatic and Syntactic" because , the compilers have made it their aim to give as much > useful information as possible concerning idioms and syntax. It is hoped that the dictionary will be of value to those who are learning English as a foreign language, not only; throughout the British Empire but also in other countries throughout the word.

Vocabulary  
The words selected for inclusion in this volume are those that the foreign student of English is likely to meet in his studies up to the time when he enters a university. After this he will probably need a dictionary with a larger vocabu- lary for reading purposes, though for writing English he will continue to find this volume useful.

Most archaic words, or those which are likely to occur only in purely scientific and technical contexts, have been excluded. Colloquial and slang words and expressions have been included if they are of the sort likely to be found in "cqiin words and Latin words and phrases of common books (e. g. modern fiction and drama) read by students. -cnce in English have also been included.

Definitions  
Dfo/initions have been made as simple as possible. Where inition in easy, common words was not practicable or .s'ltisfactory, pictures and diagrams have been supplied. A lobster, in the Cojicise Oxfoj'd Dictionary is defined as "a large marine stalk-eyed ten-footed long-tailed edible crustacean with large claws formed by the first pair of feet, bluish-black before and scarlet after boiling; it's flesh as food". I he foreign student of English, if he is a beginner, is likely to be puzzled by certain words in this definition ("stalk-eyed" and "crustacean", for example). The C. O. D. was not written specially for him. The ordinary user of the C.O.D. is likely to be a person who knows quite well what a lob- ster is and who refers to this word in a dictionary only when he needs exact and scientific information of the kind "iven in the above admirably concise and complete definition. I he foreign student usually needs only to identify the new

Correct sentence

Combined results

(b) 正解結果と統合結果の比較

図3 書籍を対象としたログ化実験 (39 フレーム)

節で述べるフレーム間認識結果統合の際にある程度は抑制される。

2.3 認識結果の統合

OCR エンジンから得られた認識結果は、フレーム間で独立である。そのため、リーディングライフ・ログで利用するには、それぞれの認識結果を統合し、装着者の読んだ文章を得る必要がある。この様子を図1に示す。このように認識結果の行対応を図りながら、結果を統合することになる。本稿では、その方法として一種の2次元 DP マッチングを用いる。

この処理の概要を図2に示す。2つの認識結果 X と Y がそれぞれ時刻 t-1 と t のフレームで得られたとする。また、x<sub>i</sub> と y<sub>j</sub> は、認識結果 X と Y 内の各文字行とする。同図の例においては、x<sub>3</sub> と y<sub>1</sub>, x<sub>4</sub> と y<sub>2</sub>, x<sub>5</sub> と y<sub>3</sub> をそれぞれ対応づけることができれば、この2フレーム間を統合できることになる。

ここで注意すべきは、OCR の出力は必ずしも常に正

解ではない点である。すなわち、X と Y も完全に正解ではなく、ある行内に文字や単語の欠落や挿入、置換が発生しうる。さらに、1行全体が欠落したり、文字に紛らわしい部分を文字としてしまうことで不要な1行が挿入されてしまう可能性もある。要するに行内と行間の両方向で(すなわち二次元的に)編集距離的な距離尺度を用いて、行間の類似度(相違度)を評価し、対応付けを行う必要がある。

そこで、図2にあるように、DP マッチングを行い、X と Y の間の最適対応関係を求める手順について述べる。文字行 x<sub>i</sub> と y<sub>j</sub> 間の距離を d(i, j) とする。この時、X と Y の最適対応付けは次の DP 漸化式を全ての i, j で計算することで求まる。

$$D(i, j) = d(i, j) + \min \begin{cases} D(i-1, j-1) \\ D(i-2, j-1) + \alpha \\ D(i-1, j-2) + \alpha \end{cases} \quad (1)$$

この漸化式計算後、(I, J) を起点として各 (i, j) で成さ

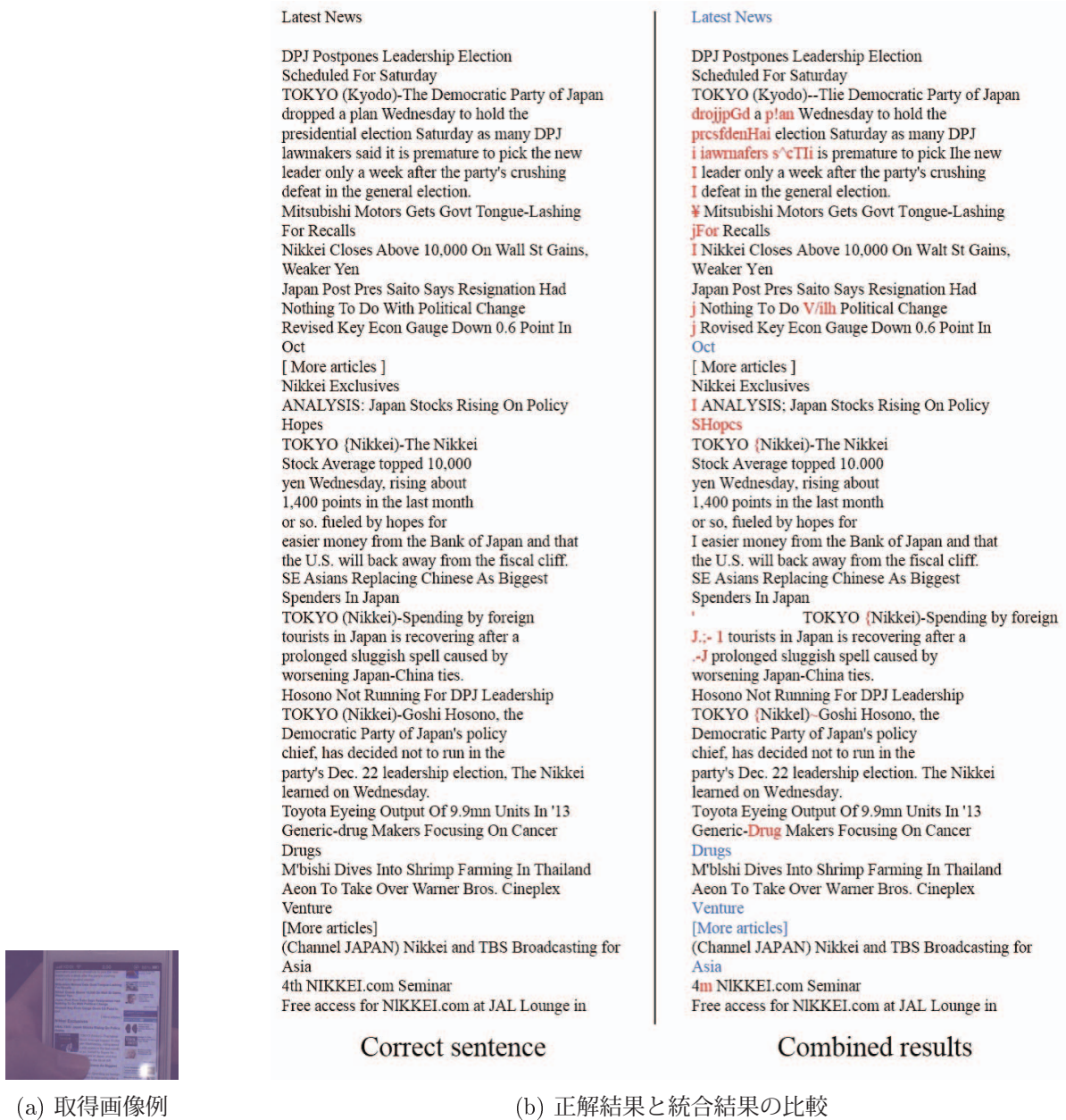


図 4 スマートフォン画面内文字を対象としたログ化実験 (29 フレーム)

れた最小値選択結果を逆戻り (バックトラック) することで、最適対応付けが求まる。なお、ここで  $\alpha$  は文字行飛越しに関するペナルティであり、文字行飛越を回避するためのものである。これにより、不要な行挿入や行欠落にロバストな対応付けが可能となる。

さらに、局所距離  $d(i, j)$  の計算には一般的な編集距離を用いる。良く知られているように編集距離も、DP マッチングによって求めることができる。これにより、行内の文字の挿入や欠落にロバストな行間距離を与えることが可能となる。以上よりわかるように、DP マッチングによる複数文字行 (X) と複数文字行 (Y) の対応付けの内部で、DP マッチングによる行単位の対応付けが行われ、編集距離  $d(i, j)$  が計算されている。すなわち二重に DP マッチングを行っていることになる。

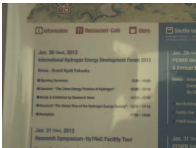
なお、DP マッチング時の端点処理には注意が必要である。すなわちカメラが文書を上から下に撮るような場合、X の上部に Y に存在しない文字行があり、Y の下部には X に存在しない行があり得る。従ってこれらをマッチングするには、いわゆる端点フリーマッチング [3] の考え方が必要となる。さらに上下でなく左右に移動することもあるので、編集距離  $d(i, j)$  の計算時にも端点フリーマッチングの考え方が必要になる。

### 3 実験

上記のシステムの動作の検証を行うため、様々な環境下で文字を認識し、認識結果の統合を行った。図 3~5 に 3 種類の結果を示す。図 3 は書籍、図 4 はスマートフォンの液晶画面、図 5 は屋外看板、である。各図に

Information Restaurant・Cafe Store  
 Jan. 30 (Wed), 2013  
 International Hydrogen Energy Development Forum 2013  
 Venue: Grand Hyatt Fukuoka  
 Opening Ceremony 9:30~10:00  
 Session1 "The Clean Energy Promise of Hydrogen" 10:00~12:10  
 Break & Exhibition by Research team 12:10~13:30  
 Session2 "The Global Rise of the Hydrogen Energy Society" 13:30~17:15  
 Reception 17:30~19:00  
 Jan. 31 (Thu), 2013  
 Research Symposium・HyTRec Facility Tour  
 Venue: Kyushu University Ito Campus  
 1) HYDROGENIUS & PCNER Research Symposium  
 international Symposium of Hydrogen Polymers Team  
 Room: Inamori Center 1F Hall 10:30~17:15  
 2013 HYDROGENIUS Tribology Symposium  
 Room: West 4 Bldg. room 312 10:45~17:40  
 Workshop on Thermal Issues for Hydrogen Energy Systems  
 Room: Inamori Center 2F Seminar room 9:45~16:20  
 Joint HYDROGENIUS and PCNER International Workshop on  
 Hydrogen-Materials Interactions  
 Room: Inamori Center 1F Hall 9:00~18:10  
 HYDROGENIUS & PCNER Joint Research Symposium  
 Fuel Cell and Hydrogen Production Symposium Jan. 28 (Mon), 2013  
 Room: Inamori Center 1F Hall 10:00-17:30  
 \*This symposium will be opened on Jan. 28th.  
 2) HyTRec Facility Tour 14:30~15:30  
 (Shuttle bus will leave 13:30 at Hydrogen Station)  
 Venue: Hydrogen Energy Test and Research Center (HyTRec)  
 \*Shuttle bus service between Kyushu University Ito Campus and HyTRec will be provided.

Information (D) Restaurant\* Cafe Store  
 Jan. 30 (Wed), 2013  
 International Hydrogen Energy Development Forum 2013  
 Venue: Grand Hyatt Fukuoka  
 Opening Ceremony 9:30~10:00  
 Session1 "The Clean Energy Promise of Hydrogen" 10:00~12:10  
 Opening Ceremony 9:30~10:00  
 Break & Exhibition by Research team 12:10~13:30  
 Session2 "The Global Rise of the Hydrogen Energy Society" 13:30~17:15  
 Reception 17:30~19:00  
 Jan. 31 (Thu), 2013  
 Research Symposium・HyTRec Facility Tour  
 Venue: Kyushu University Ito Campus  
 1) HYDROGENIUS & PCNER Research Symposium  
 international Symposium of Hydrogen Polymers Team  
 Room: Inamori Center 1F Hall 10:30~17:15  
 2013 HYDROGENIUS Tribology Symposium  
 Room: West 4 Bldg. room 312 10:45:17:40  
 Workshop on Thermal Issues for Hydrogen Energy Systems  
 Room: Inamori Center 2F Seminar room 9:45~16:20  
 Joint HYDROGENIUS and PCNER International Workshop on  
 Hydrogen-Materials Interactions  
 Room: Inamori Center 1F Hall 9:00-18:10  
 HYDROGENIUS & PCNER Joint Research Symposium  
 2013 HYDROGENIUS Tribology Symposium  
 Workshop on Thermal Issues for Hydrogen Energy Systems  
 Jan. 31 (Thu), 2013  
 Room: West 4 Bldg. room 312  
 10:45-17:40  
 Fuel Cell and Hydrogen Production Symposium Jan. 28 (Mon), 2013  
 Room: Inamori Center 1F Hall 10:00-17:30  
 \*This symposium will be opened on Jan. 28th.  
 2) HyTRec Facility Tour 14:30~15:30  
 (Shuttle bus will leave 13:30 at Hydrogen Station)  
 Venue: Hydrogen Energy Test and Research Center (HyTRec)  
 Shuttle bus service between Kyushu University Ito Campus and HyTRec will be provided.



Correct sentence

Combined results

(a) 取得画像例

(b) 正解結果と統合結果の比較

図5 屋外看板上の文字を対象としたログ化実験 (49 フレーム)

において (a) は実験対象とした動画像中の 1 フレームを示す。総フレーム数についてはキャプションに示している。フレーム画像の大きさは 2080×1552 画素であり、読む対象とした文字はおおよそ 40×50 画素程度であった。このフレーム画像は約 5fps で撮影され、視線の動きにより (スマートフォン画面については、指によるスクロール操作により)、フレームごとに何らかの方向にシフトする。

各図 (b) において、左の文字列は正解文章、右の文字列は統合結果である。図中の黒色文字は正解、赤色文字は誤認識、青色文字は認識結果が無いもの、緑色文字は前後の行の入れ替わりを表している。つまり、黒い文字が多ければ、正しい統合結果が得られたことになる。

書籍 (図 3) については、屋内環境ではありながら、カメラで撮った文面に対して、非常に高い精度で認識結果が得られていることがわかる。書籍という、比較的規則的に行配置された文書であり、さらにフォントも一般的なものであるため、市販 OCR が扱いやすい対象であったと言える。またこの結果において 40 フレームの隣接フレーム間には相当のオーバーラップがあった。(このことは、1 フレームあたりに映り込む行数が同図 (a) のように 10 行以上であることからわかる。) 一か所を除き、こうしたオーバーラップ部の冗長性はフレーム

統合の段階で適切に排除されていることもわかる。なおその一か所の失敗は、行の入れ替わりであった。

スマートフォンの液晶画面 (図 4) については、若干のノイズが見られるものの、書籍の場合と大きく変わらないほどの認識精度が得られている。指によるスクロールがあり、フレームによっては部分隠蔽が起きているが、適切なフレーム統合によりその影響もほとんど見られない。なお、各行左端に見られる “i” や “j” は画面の縁を文字と認識してしまったためである。

一方、屋外環境の看板については、大きな欠落部や誤認識の多発が確認される。これは、まず、看板上の反射により、正しく 2 値化が行えず、認識できない部分が広範囲に発生したことが挙げられる。こうした欠落があるフレームで生じると、本来「糊しろ」として機能すべき文字行が欠落してしまうことになり、二次的にフレーム間認識結果の統合にも失敗してしまうことになる。これについては文字認識結果のテキスト情報に依拠したフレーム統合を行っている以上、不可避の問題である。従って今後は、反射成分を除外するような前処理の導入はもちろのこと、トラッキングなどで画面全体の大局的移動を推定し、それをを用いてフレーム統合を制御するなどの工夫が必要である。

## 4 まとめと今後の課題

本稿では、ヘッドマウントカメラを用いたリーディングライフ・ログシステムを提案し、動作の検証実験、精度の確認を行った。その結果、屋外環境中の文字においては改良の余地が多くあることが判明するも、日常的に我々が読んでいる書籍やスマートフォン画面上の文字などは、スキャン文書をOCRしたときと大きく変わらない程度の精度で読めることが確認された。以上が確認できたことは、ヘッドマウントカメラで手持ちの文書等を撮影しても文書が認識できる、ということを実証した点で重要と考える。また、複数フレームに渡って写り込む冗長性を適切に除外できることを示せた点も成果の一つと考える。

今後の課題は、情景内文字の認識を前提とした前処理の高度化、飾り文字への対応が挙げられる。また先述の通り、画面の大局的移動を推定するためのトラッキングの利用も考えられる。さらに複数のウインドウが存在するPCディスプレイのような場合、すなわちかなり複雑なレイアウト構造をもった対象についても、一考が必要であろう。以上と並行して、ある一定のものを読んでいるから、目を離したり、書籍のページがめくられたりするような状況、謂わば「(リーディング)シーンチェン」の自動判定も、実用のためには重要な開発要素である。

## 謝辞

本研究の一部は、JST 戦略的創造研究推進事業チーム型研究 (CREST) 「共生社会に向けた人間調和型情報技術の構築」に依った。

## 参考文献

- [1] 木村崇志, 柿迫良輔, フォンヤオカイ, 内田誠一, 岩村雅一, 大町真一郎, 黄瀬浩一, “Reading-Life Log のプロトタイプ実装”, MIRU2012, 2012.
- [2] J. Liang, D. Doermann, and H. Li, “Camera-Based Analysis of Text and Documents: A Survey”, *International Journal on Document Analysis and Recognition*, vol. 7, no. 2-3, pp. 84-104, 2005.
- [3] S. Uchida, “Text Localization and Recognition in Images and Video”, in *Handbook of Document Image Processing and Recognition*, D. Doermann and K. Tombre, Eds., Springer, To be published in 2013.
- [4] 内田誠一, “DP マッチング概説 ～基本と様々な拡張～”, 信学技報, PRMU2006-166, 2006.