

るのは、この構想中のネットワーク版の SMART-GS が使われるようになったときであろう。

### 3.5 Reading unreadables : 読めないものを読む

人文学における文書研究の最大の目標の一つは、「Reading unreadables : 読めないものを読む」をいうことであろう。通常、文書とは苦もなく読めるものだと前提されている。一つの文書が何通りにも読めるというのは、その意味だけであって、それがどういう文字の並びであるかには、人によって意見の相違などない、というのが、多くの人々の常識であろう。

しかし、「テキスト分析」「テキスト理解」についての、この前提は、人文学の研究においては必ずしも成り立たない。実際、京大文学部・文学研究科が伝統的にその力を発揮し続けている分野の多くは、読めない文書を読む、そのような文献研究なのである。本論文で提案するものは、その「読めないものを読む」という、人文学においても最も困難な研究をサポートするツールなのである。

このことを理解すれば、我々の計画と類似に見える、手書き OCR や、他のチームの手書き文書サーチ研究が目指しているものが、我々の目指すものと大きく異なっていることがわかる。誤解を恐れず、一言でいえば、これらの類似研究は、万人向けのツールを目指している。一方で、我々は研究者のためのツールを目指しているのである。

そして、実現がより困難に思える「読めないものを読む研究者用のツール」の方が「素人でも使える、誰でも読める文書用の検索ツール」より現在の技術からすると、実現が遙かに容易であるという、一見矛盾しているかのようにさえ思える、この事実こそが、本プロジェクトの根幹的着想であり、また、その現実性と実効性の源なのである。

### 3.6 テキスト・イメージ・サーチと OCR

SMART-GS を紹介する際、良く受ける質問が、文字認識、つまりを使う方が良いのではないかという質問である。たとえば、折角、このようなツールを開発しても、手書き文字 OCR 等の技術が将来非常に発展して、本プロジェクトの成果が不要になるということはないだろうか？もし、そのような可能性が高いとすれば、本プロジェクトを推進する意味は低くなる。しかし、これは文献研究においては、実は本質的にありえないことなのである。その説明をしよう。

ある文書に対して OCR が存在する（既にある、開発可能である）と考えるとき、その文書について少なくとも次の二つの前提が仮定されている。

1. 一意性 : 「この OCR のこの文書に対する精度は 99%である」という言明には、その文書の読み方は一意であるということが前提されている。「正しい読み方」が存在しているから、それからの「ずれ」を読み取り精度として考えることができるのである。
  2. コード可能性 : OCR とはイメージとしての文書を文字コードの並びに変換するツールである。つまり、OCR はターゲットとなる言語のキャラクタのセットが確定して、それを表現する文字のコード体系が存在していることを前提としている。
- 一見、当然のように見える、この二つ条件のどちらも人文学研究では成り立たないことに注意しよう。

まず、一意性であるが、人文学研究では、一つのテキストの解釈が一

