

TYPE OF INDUSTRY

情報通信研究機構

NICT 先端研究

88

NICTでは自動翻訳への活用などのため、アジア言語における文字・音声処理技術に関する研究開発を行っている。私はその中でもアブギダ系文字を効率的な入力を実現す

アブギダ系文字入力効率化

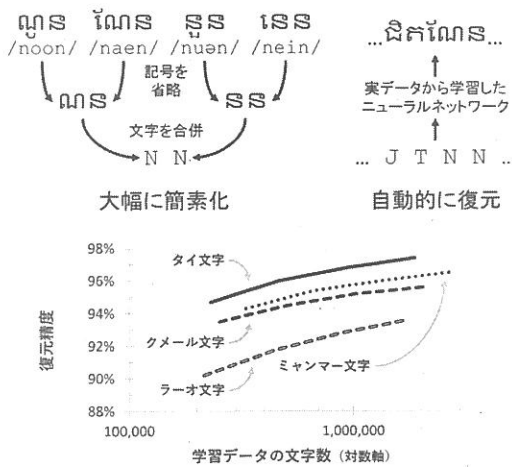
先進的音声翻訳研究開発推進センター・先進的翻訳技術研究室テニユアトラック研究員 **丁 塵辰**

15年3月筑波大学大学院修了。同年4月よりNICTに勤務、有期研究員。18年4月よりテニユアトラック研究員、現在に至る。自動翻訳、アジア言語の整備・解析に関する研究に従事。博士（工学）。



るため、計算言語学系も異なる。例えば和されることに對し、アていない。アブギダ系と記号の種類数が多い類似である独立文字を量的実データから「復論と言語処理技術を併文に使用される仮名はアジアには複雑な文字体系は「音素音節文のため、キーボード上の区別しないことにす用し、文字体系にある「音節文字」を言い、子系が用いられる。日本字」とも言う。仮名の配置が洗練されないとする。これにより元のつる。簡素化されたつづ冗長性を低減して、使音・母音ごとに記号をでは、進んだ和文入力のような「独立文字」がころも問題点となる。づりが大幅に簡素化でりに自動的・高精度のいやすい入力インター設けるローマ字は「音技術がすでに日常生活あり、母音の変換・子音が結合を表す「付加式は、付加記号を省略のニューラルネットワ復元を実現した。素文字」を言う。音節・に浸透している。音の結合を表す「付加式は、付加記号を省略のニューラルネットワ復元を実現した。一般的に文字は言語音素文字はともに「表一方、人口が多く経記号」もある。従来のし、読みが同一またはーク技術を駆使し、大アのクメール文字を例として、簡素化・復元のイメージを示している。下部はタイ文字・ミャンマー文字・クメール文字・ラオ文字における復元精度を示している。利用される実データ量の増加に従い、復元性能が向上する。現在、私たちは、実用化に向け、入力ソフトウェアを開発し、アジア言語処理技術の向上を目指している。

アブギダ系文字の簡素化および自動復元



（火曜日に掲載）

科学技術・大学