

人間と生成AIの文章理解のちがい

一年四組 高橋 尊斗

現在の生成AIは「記号接地」を全くせずに言語を学んでいるという。人間は知っている言葉が指す対象を知っており、自分自身の体から、概念の様々な特徴を捉え知ることができる（記号接地をして言葉を理解できる）。それに対して生成AIは、身体を持たず、接地した経験がないため、感覚に接地していないまま（真の意味を理解していないまま）記号と記号を結びつけ、置き換えて言葉を理解しているかのように見せてている。生成AIを活用するためには、その特性を知らなければならない。人間と生成AIとの文章理解の方法のちがいをふまえ、両者が共存する未来について考えてみたい。

二〇一八年、Google社によってBERTという文章の並び方に確率を割り当てる確率モデルである「言語モデル」が発表された。その言語モデルのうち「計算量」「データ量」「モデルパラメータ数」の三要素を大規模化したものを「大規模言語モデル」（以下、「LM」）と言い、代表的な例としOpenAI社のGPTを挙げることができる。つまり、現在のチャットボット等の生成AIは、与えられたタスクを達成するために確率の高い言葉を次々に出力しているだけであり、真にその言語を「理解」しているわけではない。

では、人間が文章を理解する際、どのような構造に注目しているのだろうか。人間は文章に一貫性を求め、他の知識を照合したり、一度構成したものから新たな問や矛盾を発見し解消したりすることで、「理解」を深めていくことができる。特に一貫性は情報間に矛盾があることに気づき、言葉や世界に関する知識（常識）である既有知識と照合することで、不自然さを吟味し解消する。³

「不自然さを解消する」ことは「矛盾の解消」と呼ばれる。情報の整合性がいつたん脅かされても適切な事例を説明としてつけ加えることによって、それらの情報を整合性あるものとして理解し直すことは、小学生でも可能であり、幼児期においてもすでに使用していることが示唆されている。⁵

この事例について、大学生を被験者として、与えられた情報の矛盾がどのように修正されるかという実験結果がある。大学生と一対一の面接を行い、例文を一度読み聞かせ、話の中につじつまの合わない所があつたか矛盾点を指摘させる。大学生は矛盾を指摘した直後にその情報解釈について「本当にそれで良いか」と念を押す。その後（イ）確かにそう言える

と思うか、（口）そう言えるとは限らないと思うが説明できないか、（ハ）そう言えるとは限らない事を説明できる（この場合に限り大学生に説明を求める）かのいずれかひとつを選択させる。最後に例文が全部本当のことだとしたらどう理解すればよいか説明させ、矛盾を解消させる。この結果、人間は最初に与えられる文の表現が変わると、矛盾の解消の際にその文の解釈を変えることが傾向として示された。

この実験と同様のものを、生成AIを使って検証してみる。使用するLLMの種類はOpenAI社のGPT-3・5、GPT-4、Google社のPaLM2、Meta社のLaMAの四種類である。なお、PaLM2は日本語での出力に応答しなかつたため、言語は英語を使用した。具体的な方法として、前項で取り上げた論文に掲載された例文を入力し、人間に對して行つたのと同じように質問を重ねて出力させた。

結果、人間でみられていた各例文の規則性が生成AIでは見られなかつた。よつてLLMに對して与える文を変えて、矛盾の解消の際にその文の解釈を変えることは難しそうである。これは、冒頭で述べたとおり、生成AIの文章理解の方法が人間と異なるからであると考えられる。また、LLMの性能が上がるほど矛盾の解消の構成数は大きく増加し、より生成AI開発の進んでいる英語で行うと構成数は増加する傾向であることも分かつた。一方、ヒントを与えることで言及カテゴリ数、回答数ともに増加する傾向があることが分かつた。暗黙裏に想定されている条件を意識化されることで、その条件が本当に満たされているのかを疑わせ、その結果、矛盾の解消促進に繋がる傾向があるといえる。

注目すべき点として、生成AIはLLMの種類によつて出力結果が変わつていたことがある。回答のあと、本当に妥当なのかと聞くと、最初は矛盾していると答えていたのに、「指摘が間違つていた」や「矛盾している点は實際にはなかつた」等、間違つた主張に変えてしまつ出力パターンが多く確認された。最初は矛盾していないと回答していたLLMが、あとになつて矛盾していたと訂正するのである。

以上より、生成AIは本質的な「理解」をしていてるかのように見せる状態にはほど遠いことがわかる。しかし、一部の側面では人間と比べて、文章理解の能力が同等もしくはそれ以上と見なしうる場合もあつた。年々、人間とは異なつた仕方で「理解」する、しているように見せる能力が向上しているため、数年後、AIの「理解」が、人間のその能力を上回るようになるかもしれない。

なお、現在開発途中ではあるが、「なぜその答えを出したのか」を説明できるAIは、「説

明可能なAI（Explainable AI）」と呼ばれている。説明可能なAIという概念には、AIの出す結論が信頼に足るものなのか最終的に判断し、責任を持つのは人間であるという前提がある。しかし、AIの発展によってAI自身が「なぜその答えを出したのか」を説明し、AIモデルの説明可能性を自身が証明するようになるかもしれない。それが実現可能にすれば、人間が介さずとも自身を説明し、自己認識をしているかのように見せることができるということだ。

今後もAIは加速度的に発展し続けるだろう。電話交換手のように、新しい道具、技術の発展によって消失する職種は多数存在する。自身が「理解」する能力を超えた人は、現時点でも仕事を奪われることは確実だ。現在のAIは人間による入力が行われないと動かず出力しないが、自分で自我を持つかのような説明可能なAI、汎用AIが生まれた後には、真のシンギュラリティが待っているかもしれない。

- 1 今井むつみ・秋田喜美『言語の本質』(1999年・中央公論新書)
- 2 NRI「大規模言語モデル」<https://www.nri.com/jp/knowledge/glossary/lst/ta/lm> (1999年7月二十二日閲覧)
- 3 内田伸子「III・文章理解と知識」(佐伯胖『認知心理学講座第3巻推論と理解』1982年・東京大学出版会)
- 4・6 鈴木孝子「一見矛盾する情報の理解過程における事例構成」(『教育心理学研究 33巻2号』1985年)
- 5 久保ゆかり「幼児における矛盾する出来事のエピソードの構成による理解」(『教育心理学研究 30巻3号』(1982年))