

生成 AI のトークンと形態素解析の関係

◇ トークンとは？

生成 AI における「トークン」は、モデルがテキストを処理する最小単位です。これは必ずしも「単語」ではなく、以下のような単位になることがあります：

- 単語（例：apple）
- 単語の一部（例：un, believ, able）
- 記号や空白（例：. や ）

◇ 形態素解析とは？

形態素解析は、自然言語処理（NLP）において、文を意味のある最小単位（形態素）に分解する処理です。日本語のように単語の区切りが明示されていない言語では特に重要です。

例：

「私は学生です」 → 「私 / は / 学生 / です」

形態素解析では、品詞（名詞、動詞、助詞など）も同時に判定されることが多いです。

◇ 両者の関係

- 共通点：どちらも「テキストを小さな単位に分割する」処理です。
- 違い：
 - トークン化はモデル内部の処理に最適化されており、意味的な区切りとは限りません。
 - 形態素解析は言語学的な意味に基づいて分割されます。

◇ 生成 AI と形態素解析の関係

- 日本語などの言語では、形態素解析を使って前処理を行い、その結果をトークン化することがあります。
- ただし、最近の大規模言語モデル（例：GPT、BERT など）は、形態素解析を使わずに直接トークン化することが多いです。

◇ まとめ

項目	トークン	形態素解析
単位	モデル内部の処理単位	意味のある言語単位
用途	モデルの学習・生成	言語解析・前処理
関係	間接的に関連あり	モデルによっては併用されることも