

# ジャンル別の代表的なコーパス

- 書籍

- ❖ Book Courpus: 様々なジャンル/トピックの11000冊以上の書籍から構成
- ❖ Project Gutenberg: 小説・エッセイ・詩・演劇・歴史・科学・哲学など70000以上の書籍から構成

- **CommonCrawl** (ペタバイトケールのデータ量を含む最大級のウェブクロールデータベース)

- ❖ C4: CommonCrawlにクリーニングを行った800GBのデータ
- ❖ RealNews: CommonCrawlからニュースを抽出した120GBのデータ

- **Reddit**

- ❖ OpenWebText: Reditで高く評価された投稿を集めた38GBのデータ

# ジャンル別の代表的なコーパス

- **Wikipedia**

- ❖ Wikipedia: 英語版の記事を集めた21GBのデータ

- **コード**

- ❖ BigQuery: 様々なプログラミング言語のコードスニペットを含むデータ

- **その他**

- ❖ Pile: 書籍・WEB・コード・論文・SNSなど多様なソースから構成される800GBのデータ