

シンポジウム

「進化する Web コーパス：現状と課題」

投野由紀夫

まえがき

本シンポジウムは、近年急速に研究領域として拡大しつつある Web コーパスの動向を紹介し、英語コーパス研究におけるその位置づけ、現状と課題を検討するという趣旨で、2010年4月24日、英語コーパス学会第35回大会（会場：兵庫県立大学）において開催された。

Web コーパスはこの10年ほどで圧倒的な速度で増え続けている。2006年9月には Google が1兆語の Web データをもとにした n-gram 頻度リスト、Google 1T 5-gram Version 1 を Linguistic Data Consortium から配布したのは記憶に新しい。Web コーパスを利用した代表的な検索システムが Sketch Engine (<http://www.sketchengine.co.uk>) である。2010年10月1日現在、Sketch Engine のトップ画面からは52種類の多言語コーパスが検索可能であるが、そのうちの約20種類が Web 上のテキストをソースとして構築されたコーパスであり、特に1億語以上の大規模コーパスはほとんどが Web コーパスである。また Sketch Engine 以外にも、一般利用者が Web コーパスを簡単に作成できるツール群（例：BootCat¹）や、Web 全体をコーパスとして検索するような検索エンジンとコンコーダンスの融合ツール（WebCorp² など）も4,5年前よりも格段に品質向上してきており、Web コーパスを取り巻く利用環境は加速度的に充実してきている。

本シンポジウムでは、講師として、田中省作（立命館大学）、大羽良（中央大学）、中村隆弘（ネットアドバンス）、星野守（小学館）の4氏をお招きし、投野がモデレーターおよび問題提起の発表を行った。まず投野が「Web コーパス概観」と題する導入的発表を行った後、田中講師が「Web コーパスの言語情報処理基盤」という題で、Web コーパス構築に関する自然言語処理の基礎技術の紹介および研究基盤としての Web コーパス、個別研究における Web コーパスの位置づけを論じた。続いて、大羽講師が「Web コーパス研究におけるブログの可能性とその文体的特徴」というテーマで、Web 上に存在するブログのもつコーパスとしての潜在的な可能性とその話し言葉・書き言葉の両面を併せ持つ文体的な特徴について実証データに基づき論じた。中村・星野両講師は「小学館 Web コーパス [Sekai Corpus] 構築とその活用」と題して、小学

館コーパス・ネットワークの一環として構築した 20 億語規模の Sekai Corpus の概要と、その専門分野英語比較の作業から得た活用可能性を論じていただいた。最後に再び投野による総括的な発表の後、パネリストおよびフロアとのディスカッションに入った。

以下に掲載した 3 編の論文は、田中氏、大羽氏、および投野の発表内容をもとに加筆修正したものである。中村・星野両氏の発表に関しては、残念ながら今回論文の形で掲載することはできなかったが、最後のディスカッションの内容などはできるだけ投野の論文に盛り込む形でまとめさせていただいた。

本シンポジウムの 3 編の総括論文が、英語コーパス研究における Web コーパスの利活用に少しでも参考になれば望外の幸せである。

注

¹ <http://bootcat.sslmit.unibo.it/>

² <http://www.webcorp.org.uk/>