# Top-down wisdom versus bottom-up noise: can data (ever) replace linguistic expertise?

**Dagmar Divjak & Petar Milin**

**University of Birmingham, UK**

Recent achievements in language engineering have challenged many of the assumptions that linguists have long taken for granted. In particular, state-of-the-art Large Language Models (LLMs) routinely succeed in a wide range of downstream cognitive tasks without relying on abstract categories, discrete levels of representation, or even exemplar-based models of usage (Linzen & Baroni, 2021). From a linguistic perspective, the success of LLMs built on deep neural networks is especially striking because their architectures and input data are neither rooted in symbolic representations nor transparently linked to real-world instances of language use. Instead, language is re-represented in the form of high-dimensional numerical vectors—mathematically elegant, but profoundly opaque to humans.

Jelinek's (1988) now-famous quip—"Every time I fire a linguist, the performance of the speech recognizer goes up"—continues to resonate. It captures a deep and ongoing tension between data-driven and theory-driven approaches to language, a tension that has only intensified as so-called 'foundation models' scale ever upwards. Nearly four decades later, this remains what Sutton (2019) calls "a bitter lesson": that brute-force learning on massive data—and pure trial-and-error compute—may outperform carefully constructed, theoretically motivated systems.

This keynote reflects on what this lesson means for linguistics. Can linguistic theory still contribute to models of language use, learning, and generalisation—or has it been superseded by scale and statistics? To what extent do our traditional abstractions help—or hinder—our ability to model language as it is actually encountered, processed, and produced? And how might linguists productively engage with systems that appear to succeed without any obvious linguistic insight at all?

Rather than offering a binary choice between symbolic expertise and statistical scale, we will argue for a reframing of the relationship between them. Using examples from our work which merges theoretical linguistics with experimental psychology and artificial intelligence, we explore how linguistic knowledge can inform the design, interpretation, and critical evaluation of data-driven models—and where it may need to evolve. In doing so, we aim to reopen the conversation between linguistics and language technology, asking not only what linguists can learn from machine learning, but what machine learning still has to learn from linguists.

## References

Jelinek, F. (1988). *Applying Information Theoretic Methods: Evaluation of Grammar Quality* Workshop on Evaluation of NLP Systems, Wayne PA.

Linzen, T., & Baroni, M. (2021). Syntactic structure from deep learning. *Annual Review of Linguistics, 7*(1), 195-212. https://doi.org/10.1146/annurev-linguistics-032020-051035

Sutton, R. (2019). The bitter lesson. *Incomplete Ideas (blog), 13*(1), 38. [http://www.incompleteideas.net/IncIdeas/BitterLesson.html]