# Using ECCO-BERT and the Historical Thesaurus of English to Explore Concepts and Agency in Historical Writing

## Interpreting the Eighteenth-century Luxury Debate

**Liimatta, Aatu**
aatu.liimatta@helsinki.fi
University of Helsinki

**Mäkelä, Eetu**
eetu.makela@helsinki.fi
University of Helsinki

**Ginter, Filip**
figint@utu.fi
University of Turku

**Rastas, Iiro**
iitara@utu.fi
University of Turku

**Tihonen, Iiro**
iiro.tihonen@helsinki.fi
University of Helsinki

**Zhang, Jinbin**
jinbin.zhang@aalto.fi
Aalto University

**Pivovarova, Lidia**
lidia.pivovarova@helsinki.fi
University of Helsinki

**Tolonen, Mikko**
mikko.tolonen@helsinki.fi
University of Helsinki

**K, Milja**
kaledriina@gmail.com
University of Helsinki

**Babbar, Rohit**
rohit.babbar@aalto.fi
Aalto University

**Wang, Ruilin**
ruilin.wang@helsinki.fi
University of Helsinki

**Säily, Tanja**
tanja.saily@helsinki.fi
University of Helsinki

**Ryan, Yann Ciarán**
yann.ryan@helsinki.fi
University of Helsinki

## Introduction

This paper introduces a new method to determine meanings in historical texts. It relies on a context-aware language model, using Transformer architecture. As a case study, we use the method to understand more about the eighteenth-century 'luxury debate'. At this time, the growth and acceptance of commercial society and an increase of the availability of luxury goods from overseas meant that the dominant discourse shifted from seeing luxury as a corrupt force to something more complex. Using this new method, we focus on two very specific meanings of luxury and their change over time: luxury as a 'disease' and as something 'productive'. The results show that the average value of luxury in this latter category increased over the century, suggesting a change in the way authors wrote about the concept.

## Related work

This method builds on existing computational approaches to understanding semantic meaning in texts, but updates them with the use of a trained-from-scratch BERT Transformer model. Many existing approaches to semantic representations focus the distributional hypothesis, which compares the contexts and distributions of words to determine semantic similarities (de Bolla et. al. 2020). However, these general models have been criticised for lacking the specificity needed by historians (Wevers & Koolen, 2020). Particularly relevant to this paper is the work done to disambiguate word senses to a set of higher-level classes or categories, usually derived from large semantic dictionaries (Ciaramit et al 2003, Piao et. al 2017; Robinson et. al 2022). BERT transformer models have also previously been used to detect shifts in historical language meaning (Martinc et al., 2020, Montariol et al., 2021).

Data and Methods

The paper primarily uses two resources: the ECCO-Bert language model created by the TurkuNLP group and the Historical Thesaurus of English (HTE). ECCO-Bert is a dedicated BERT model trained from scratch on a large collection of mostly-English-language texts, Eighteenth Century Collections Online (Rastas et. al. 2022). The second resource, the HTE, is a thesaurus compiled over four decades by experts, derived from the Oxford English Dictionary (Kay et. al, 2023). It has a detailed semantic structure, meaning each word in the thesaurus is contained within a hierarchical set of categories, each representing a more fine-grained group of connected words. Each word in the thesaurus is assigned a first and last date of known usage.

As Robinson et. al (2022), we assume that a given semantic category can be referred to by the lexical items contained within it. The concept of 'ill health,' for example, might be represented by the words *sickness*, *plague*, *ailment*, and so forth. The proposed method extracts senses from the representations of words as found in ECCO-BERT, using references from the HTE. This is done by removing a 'seed' word of interest from a sequence of text, and using the model to generate the most likely replacements, along with probability scores. For a given word, first the top *k* replacement words were generated, using the BERT model masked language task, resulting in a representation of that word. In the next step, this representation was compared to the lemmas representing each semantic category, filtered to only include lemmas in the HTE from within the set of target dates (1700–1800). The confidence score for each category is the sum of the model probabilities of a subset of the words from the representation which are found within the subcategory lemmas.

# Results

As a preliminary study, this method was applied to each instance of a single word from the ECCO dataset, *luxury*. Luxury in the eighteenth century, wrote Berg and Eger (2003) "[...] gradually lost its former associations with corruption and vice, and came to include production, trade and the civilising impact of superfluous commodities." Tracing the key semantic categories for instances of the word *luxury* allows us to better understand this changing attitude over time. Using the method, we extracted category scores for all instances of the word luxury in the ECCO dataset. In many cases, the method was able to derive sensible and meaningful semantic categories for these words. This is most apparent in the aggregate, as seen in the results from the top ten below (Table 1).

Table 1: Top ten semantic categories for the word luxury in ECCO. The score given is the sum of the confidence scores for each relevant word as given by the model.

| Semantic Category | Sum of Confidence Scores |
| --- | --- |
| the world -> physical -> sensation-> physical sensibility-> sensuous pleasure | 3677.288 |
| the mind-> emotion-> pride-> ostentation | 3463.470 |
| the world -> action or operation -> advantage -> usefulness | 3226.382 |
| society -> leisure -> entertainment -> pastimes | 3186.703 |
| the mind -> attention and judgement -> judgement or decision -> quality of being good | 3185.732 |
| the mind -> attention and judgement -> beauty -> splendour | 2906.491 |
| the world -> physical sensation -> sexual relations -> sexual activity | 2818.898 |
| the mind -> attention and judgement -> esteem -> reputation | 2701.080 |
| the world -> health and disease -> ill health -> a disease | 2685.506 |
| society -> society and the community -> social class -> nobility | 2573.009 |

Column N in this table is a sum of the confidence scores produced by the method explained above, essentially ranking the top semantic categories for all instances of the word luxury in our dataset. Of particular relevance are two from this list: **the world -> action or operation -> advantage -> usefulness** and **the world -> health and disease -> ill health -> a disease**.

These two categories are good proxies for the two dominant and opposing views of luxury: something corrupting and responsible for moral decline; or something productive and useful. Examples of high-scoring instances are given in Figure 1 below.



love or labour, which were the means by which they had obtained it. But, alas, it was not long ere Cyrus himself sowed the first seeds of that luxury which soon overspread and corrupted the whole nation : for being to shew himself on a particular occasion to his new conquered subjects, he thought proper, in order to heighten the splendor of his regal dignity, to make a pompous display of all the magnificence and shew, that could be contrived to dazzle the eyes of the people. Among other things he changed his own apparel, as also that of his officers, giving them all garments

THE WORLD. 101

The natural consequences of security and affluence in any country, is a love of pleasure ; when the wants of nature are supplied, we seek after the conveniencies; when possessed of these, we desire the luxuries of life; and when every luxury is provided, it is then ambition takes up the man, and leaves him still something to wish for : the inhabitants of the country, from primi-

*Figure 1: examples of high-scoring instances of the word luxury in ECCO texts in the semantic category 'usefulness' (left) and 'disease' (right). Left is from a letter by Oliver Goldsmith, who wrote extensively on the upsides of luxury, and right is from a work called Beauties of History, taken from a section devoted to luxury and its ills.*

To understand more about how these views of luxury changed over time, we graphed the average value of each for each year of the century (Figure 2).
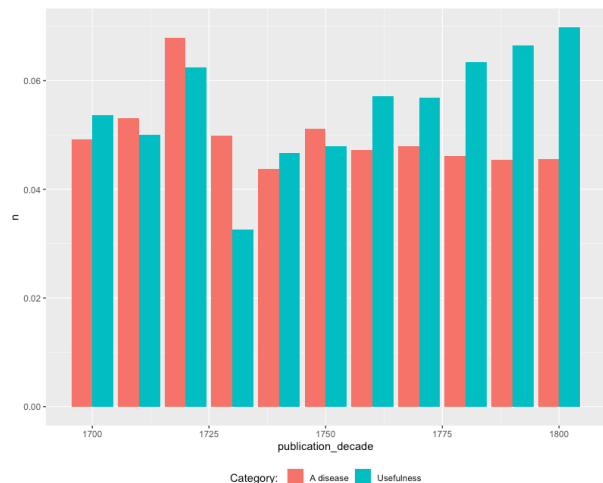


*Figure 2: Average value of the semantic category for the word luxury, across the eighteenth century.*

Though noisy, this points to a fairly flat average use of the word *luxury* within the 'disease' semantic category and an increase over time of the 'usefulness' category. This is an interesting finding and shows that there is a real statistical background to the idea that the prevailing view of luxury became less negative over time.

The new usage is not just a general 'positive vs negative', but a shift towards a very specific semantic meaning of something being productive and useful.

# Conclusion

In conclusion, this paper proposes a method for determining layered meanings in a historically-aware manner using large language models and a widely-used semantic classification system. The proposed method is applied to the word *luxury* found in a dataset of over 30 million pages of early modern English to trace the evolution of its meaning over time and the differences in use between authors. This combination of the statistical, context-aware Transformer model, and the expert semantic resource of the HTE shows great promise, and in future work we aim to use the method over a wide range of concepts and semantically-varying terms.

# Bibliography

**Berg, M., & Eger, E.** (2003). The rise and fall of the luxury debates. In M. Berg & E. Eger (Eds.), *Luxury in the Eighteenth Century* (pp. 7–27). Palgrave Macmillan UK. https://doi.org/10.1057/9780230508279_2

**Bolla, P. D., Jones, E., Nulty, P., Recchia, G., & Regan, J.** (2020). The idea of liberty, 1600–1800: A distributional concept analysis. *Journal of the History of Ideas*, *81*(3), 381–406. https://doi.org/10.1353/jhi.2020.0023

**Ciaramita, M., & Johnson, M.** (2003). Supersense tagging of unknown nouns in wordnet. *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, 168–175. https://aclanthology.org/W03-1022

**de Bolla, P.** (2020). The evolution of aesthetic concepts 1700–1800. In K. Axelsson, C. Flodin, & M. Pirholt (Eds.), *Beyond autonomy in eighteenth-century British and German aesthetics*. Routledge.

**Martinc, M., Kralj Novak, P., & Pollak, S.** (2020). Leveraging contextual embeddings for detecting diachronic semantic shift. *Proceedings of the Twelfth Language Resources and Evaluation Conference*, 4811–4819. https://aclanthology.org/2020.lrec-1.592

**Montariol, S., Martinc, M., & Pivovarova, L.** (2021). Scalable and interpretable semantic change detection. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4642–4652. https://doi.org/10.18653/v1/2021.naacl-main.369

**Piao, S., Dallachy, F., Baron, A., Demmen, J., Wattam, S., Durkin, P., McCracken, J., Rayson, P., & Alexander, M.** (2017). A time-sensitive historical thesaurus-based semantic tagger for deep semantic annotation. *Computer Speech & Language*, *46*, 113–135. https://doi.org/10.1016/j.csl.2017.04.010

**Rastas, I., Ciarán Ryan, Y., Tiihonen, I., Qaraei, M., Repo, L., Babbar, R., Mäkelä, E., Tolonen, M., & Ginter, F.** (2022). Explainable publication year prediction of eighteenth century texts with the bert model. *Proceedings of the 3rd Workshop on Computational Approaches to Historical Language Change*, 68–77. https://doi.org/10.18653/v1/2022.lchange-1.7

**Robinson, J. A., & Weeds, J.** (2022). Cognitive sociolinguistic variation in the old bailey voices corpus: The case for a new concept -led framework. *Transactions of the Philological Society*, *120*(3), 399–426. https://doi.org/10.1111/1467-968X.12250

**Wevers, M., & Koolen, M.** (2020). Digital begriffsgeschichte: Tracing semantic change using word embeddings. *Historical Methods: A Journal of Quantitative and Interdisciplinary History*, *53*(4), 226–243. https://doi.org/10.1080/01615440.2020.1760157