



A CORPUS DRIVEN COMPARISON OF ENGLISH AND UZBEK LITERARY PROSE

Satibaldieva Nigora

Teacher, Uzbek State World Languages University

Abstract

This article analyses the dynamics of lexical change in English and Uzbek literary language between the nineteenth and twentieth centuries. Drawing on a 7-million-word diachronic subset of the Corpus of Historical American English and a newly compiled 3-million-word Uzbek Fiction Corpus, we trace the emergence, decline, and semantic drift of high-frequency lemmas. Statistical measures – log-likelihood, keyness, and adjusted type–token ratio – reveal parallel patterns of technological expansion and the attrition of archaic lexis, but also divergent trajectories shaped by colonial contact, script reform, and nation-building. The findings refine existing models of lexico-semantic evolution and offer practical implications for historical lexicography and literary translation.

Keywords: Lexical change, diachronic corpus, English literature, Uzbek literature, nineteenth century, twentieth century, keyness analysis, historical lexicography

Introduction

Lexical change in literary prose provides a barometer of cultural and technological transformation. While English has served as a touchstone for diachronic corpus studies (Kilgarriff, 2005), far less is known about Turkic languages such as Uzbek, whose literary canon underwent radical shifts during Soviet modernisation (Yakubov, 2021). A bilingual comparison illuminates both universal and locally specific mechanisms of change, enriching theories of language evolution in contact zones (Thomason, 2001).

Early diachronic corpus studies in English—most notably the *ARCHER* project (Biber & Finegan, 1997) and the *Lancaster-Oslo/Bergen Corpus* expansions



(Leech, 2004)—demonstrated that lexical change is intrinsically linked to shifting register preferences and socio-cultural context. Subsequent research has refined methodological rigour: Davies and Fuchs (2015) illustrate how frequency-adjusted keyness metrics capture “slow-burn” lexical diffusion, while Hilpert (2020) advocates vector-space models to trace fine-grained semantic drift. Uzbek, by contrast, entered large-scale corpus linguistics relatively late. Yakubov (2021) provides the first systematic overview of Soviet-era lexical engineering, but his manual sampling limits quantitative generalisation. Recent projects such as the *TilKom* Diachronic Corpus (Musurmonova, 2022) and Mustafayev’s (2023) Fiction Corpus leverage OCR-corrected scans and lemmatisation tools adapted to both Cyrillic and Latin scripts, enabling token-level comparisons across orthographic regimes. Their findings confirm a post-1991 surge in Turkic-root neologisms, partly reversing decades of Russian lexical dominance.

Cross-linguistic comparisons remain scarce. Heine and Kuteva’s (2005) contact-induced change model predicts symmetrical borrowing where sociopolitical power is balanced, yet Uzbek demonstrates *asymmetrical* pathways: Russian loans proliferated in technical fields, while Uzbek preserved native vocabulary in agrarian semantics (Mukhamedov, 1999). This observation aligns with Pagel’s (2017) evolutionary linguistics claim that basic-vocabulary change slows under strong identity pressure.

From a methodological perspective, Kilgarriff’s (2005) critique of “lexical random drift” cautions against attributing statistical noise to cultural causes. Brezina (2020) recommends multi-metric triangulation—MATTR, VOCD, and log-likelihood—when analysing corpora of unequal sizes, a practice we adopt here. Finally, Kutuzov and Øvrelid’s (2020) diachronic embeddings offer a reproducible framework for measuring semantic displacement, inspiring the vector-space component of our study.

Together, these strands establish both a theoretical foundation and a set of best practices that guide our bilingual investigation of nineteenth- and twentieth-century literary lexis.



Corpus and Methodology

The English dataset comprises the fiction section of the Corpus of Historical American English (COHA), 1810-1910 and 1911-2000, balanced for genre and author gender (Davies, 2010). The Uzbek dataset draws on digitised first editions of canonical authors – Cho‘lpon, Fitrat, Aytmatov – as well as contemporary prose, each decade represented by $\approx 150\,000$ tokens. Texts were lemmatised with TreeTagger (Schmid, 1994) and, for Uzbek, manually validated to offset script-conversion noise (Mustafayev, 2023). Keyness was calculated with log-likelihood ($p < 0.001$), and lexical diversity with the moving-average type–token ratio (Brezina, 2020). Semantic shift candidates were flagged via vector-space displacement scores using diachronic word embeddings (Kutuzov & Øvrelid, 2020).

Results

Both languages exhibit a surge in technonyms after 1880. English shows significant keyness for *engine*, *telephone*, *motor*, *radio*, while Uzbek registers *zavod* “factory”, *elektr* “electricity”, and *samolyot* “airplane”, all Russian loanwords (Mukhamedov, 1999). Relative frequencies rose 450 % in English and 390 % in Uzbek between Period 1 (1810-1910) and Period 2 (1911-2000), mirroring industrialisation timelines.

Obsolete English markers such as *thou*, *hath*, *ere* fell below 2 instances per million words by the 1950s. Uzbek displayed a comparable decline in Persian-Arabic archaisms like *hushyor* “alert” and *anjuman* “assembly”, dropping 78 % after 1930, coinciding with Soviet lexical purification policies (Yakubov, 2021).

Vector-space analysis pinpointed forty English lemmas with high cosine displacement; seminal cases include *broadcast* (from “sowing by scattering” to radio/TV “transmission”) and *gay* (from “joyful” to “homosexual”). Uzbek highlighted *komsomol* (originally “Communist youth union”) now metonymic for nostalgia, and *qo‘ltelefon* shifting from “landline handset” to “mobile phone”. Drift trajectories correlate with socio-political realignments and technological adoption lags (Mustafayev, 2023).

English fiction maintained a stable MATTR ($0.73 \rightarrow 0.72$), suggesting lexical replacement rather than inflation. Uzbek MATTR increased from 0.69 to 0.75,



reflecting influxes of Russian-derived neologisms and later, post-independence revival of Turkic roots.

Parallel technological growth supports Labovian models of lexical diffusion triggered by referential need (Labov, 2001). Yet the pathways diverge once sociolinguistic variables intervene: Russian's mediator role yielded "borrowed modernity" in Uzbek, whereas English generated neologisms internally. Script reform (Arabic → Latin → Cyrillic → Latin) further accelerated Uzbek lexical turnover, a variable absent in the English timeline. These findings corroborate Heine and Kuteva's (2005) contact-induced change theory, extending it to macro-lexical scale.

A diachronic corpus lens uncovers both convergent and divergent patterns in English and Uzbek literary lexis across two centuries. Shared technological expansion contrasts with language-specific forces – imperial contact, orthographic upheaval, nation-building – that shape lexical retention and semantic drift. Future research should integrate genre-specific subcorpora (drama, memoir) and align morphological tagging to refine cross-language comparability.

References

1. Brezina, V. (2020). *Statistics in corpus linguistics: A practical guide*. Cambridge University Press.
2. Davies, M. (2010). The Corpus of Historical American English: 400 million words, 1810-2009. *International Journal of Corpus Linguistics*, 15(2), 301-310.
- Heine, B., & Kuteva, T. (2005). *Language contact and grammatical change*. Cambridge University Press.
3. Kilgarriff, A. (2005). Language is never, ever, ever, random. *Corpus Linguistics and Linguistic Theory*, 1(2), 263-276.
4. Kutuzov, A., & Øvrelid, L. (2020). Temporal referentiality in diachronic word embeddings. *Digital Scholarship in the Humanities*, 35(suppl. 1), 42-58.
5. Labov, W. (2001). *Principles of linguistic change, Vol. 2: Social factors*. Blackwell.



6. Mukhamedov, U. (1999). Russian lexical influence on modern Uzbek prose. *Philologia*, 3, 45-57.
7. Mustafayev, B. (2023). Building a diachronic Uzbek fiction corpus: Challenges and solutions. *Language Resources & Evaluation*, 57(2), 505-526.
8. Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. *Proceedings of the International Conference on New Methods in Language Processing* (pp. 44-49).
9. Thomason, S. (2001). *Language contact: An introduction*. Edinburgh University Press.
10. Yakubov, J. (2021). Lexical standardisation in Uzbek: From Soviet modernisation to post-independence revival. *Journal of Central Asian Linguistics*, 12(1), 23-41*.
11. Dalieva M. K. et al. The function of english songs to improve listening skill //NovaInfo. Ru. – 2021. – №. 124. – С. 30-32.
12. Dalieva M. INTEROPERATION OF LANGUAGE, SCIENTIFIC TERMINOLOGY, AND INTERDISCIPLINARY COLLABORATION //Western European Journal of Linguistics and Education. – 2024. – Т. 2. – №. 1. – С. 1-4.
13. Dalieva M. K. et al. The theories of teaching vocabulary in context //NovaInfo. Ru. – 2021. – №. 124. – С. 45-46.
14. Satibaldiev E., Kuldashv A. LEXICAL PROBLEMS OF TRANSLATING THE NAMES OF NEW PHENOMENA IN SOCIAL LIFE //Студенческий форум. – 2020. – №. 19-3. – С. 86-88.
15. Dalieva M. X., Satibaldiev E. K. WAYS OF ELIMINATING POLYSEMY IN THE LANGUAGES OF DIFFERENT SYSTEMS //ББК 81.2 я43 Методика преподавания иностранных языков и РКИ: традиции и инновации: сборник научных трудов VIII Международной научно-методической онлайн-конференции, посвященной Году педагога и наставника в России и Году русского языка в странах СНГ (11 апреля 2023 г.)–Курск: Изд-во КГМУ, 2023.–521 с. – 2023. – С. 35.
16. Сатибалдиев Э. К. ЯЗЫКОВОЕ КОНТАКТИРОВАНИЕ: БИЛИНГВИЗМ, ПОЛИЛИНГВИЗМ, ИНТЕРФЕРЕНЦИЯ



-
- //ИНОСТРАННЫЙ ЯЗЫК В ПРОФЕССИОНАЛЬНОЙ СФЕРЕ:
ПЕДАГОГИКА, ЛИНГВИСТИКА, МЕЖКУЛЬТУРНАЯ
КОММУНИКАЦИЯ. – 2022. – С. 144-149.
17. Далиева, М. (2024). ОСОБЕННОСТИ ПОЛИСЕМИИ КАК
КОНЦЕПТУАЛЬНОГО ФЕНОМЕНА. TAMADDUN NURI JURNALI,
5(56), 508-510.
 18. Сатибалдиев Э. К. Родной и неродной языки: лингвистические и
методические аспекты речевой интерференции. – 2022.
 19. Сатибалдиев Э. К. Semantic Change and Its Sources in the English
Language //Филологические науки в России и за рубежом. – 2019. – С.
15-16.
 20. Далиева, М. (2024). Comparative and typological approaches to analyzing
polysemy in linguistic terms. Актуальные вопросы языковой подготовки
в глобализирующемся мире, 1(1).
 21. Далиева, М. (2023). POLYSEMY IN COGNITIVE LINGUISTICS.
American Journal of Pedagogical and Educational Research, 10, 138–140.
 22. Далиева, М. (2024). Когнитивные модели полисемии лингвистических
терминов. Каталог монографий, 1(1), 1-153.
 23. Satibaldieva N. Polysemy of terms in computational linguistics //International
Journal of Scientific Trends. – 2024. – Т. 3. – №. 1. – С. 82-84.
 24. Сатибалдиева Н. The function of neologisms in language progression
//Актуальные вопросы языковой подготовки в глобализирующемся
мире. – 2024. – Т. 1. – №. 1.
 25. Сатибалдиева Н. The effectiveness of parallel and diachronic corpora in
modern language-related research //Лингвоспектр. – 2024. – Т. 4. – №. 1. –
С. 9-17.