

Дарчук Н. П. Возможности семантической разметки корпуса украинского языка (КУЯ).

В статье рассмотрены лингвистические основы семантической разметки Корпуса украинского языка как четвертого этапа представления информации о единицах Корпуса. В основу разметки положена таксономическая классификация корпуса русского языка, но дополненная и видоизмененная. Создано программное обеспечение для работы в он-лайн режиме. Материалом послужил частотный словарь публицистического стиля объемом в 40 тыс. лексем, созданный на выборке в 16 млн словоформ украиноязычного текста.

Ключевые слова: Корпус текстов, семантическая разметка, таксономическая классификация, таксон.

Darchuk N. P. Capabilities of Semantic Tagging Within the Ukrainian Corpus.

The article views linguistic aspects of semantic tagging within the Ukrainian Corpus. The lexical content of texts of different genres, in particular, modern fiction, drama, journalism, scientific, popular scientific, and business will be provided with a specific tagging respectively. The work represents two types of tagging: I – a taxonomic one, featuring journalistic and fiction genre and II – a thesaurus-based tagging specifically for scientific and business genres.

The tagging is based on taxonomic classification applied in the Russian Corpus but extended and extra modified. There were developed the software tools for online work based on materials of frequency dictionary of journalistic style with a total volume of 40,000 lexems compiled from the sampling of 16 Million word forms of Ukrainian texts. The thesaurus-based approach is grounded on the identification of thematically relevant lexical-semantic variations and grouping them by applying a formalized method of a thesaurus construction, which meets the standards of modern terminography. There were developed the software tools for performing of two types of semantic tagging.

Keywords: linguistic corpus, semantic tagging, taxonomic classification, taxon, thesaurus, information retrieval system.

УДК 81'33(477)

Дарчук Н. П., Лангенбах М. О., Сорокін В. М., Ходаківська Я. В.
Київський національний університет
імені Тараса Шевченка

ПАРАЛЕЛЬНИЙ КОРПУС ТЕКСТІВ ПАРКУМ

Незважаючи на активний розвиток корпусної лінгвістики в Україні, досі існує велика прогалина в царині розробки паралельних корпусів. Метою роботи є формулювання основних засад творення та використання паралельного корпусу текстів ПарКУМ. Завдання, що вирішуються в ході дослідження: визначення напрямів перекладу та принципів добору текстів; вибір основних параметрів розмітки; визначення концепції роботи з матеріалом; розробка структури користувацького інтерфейсу. Подається інформація про всі типи розмітки, передбачені в корпусі: метатекстову, структурну та лінгвістичну. Окрім опису структури проекту, подано роз'яснення щодо принципів роботи з корпусом.

Ключові слова: паралельний корпус, корпусна лінгвістика, корпусна розмітка, паралельні тексти, перекладні відповідники.

Серед сучасних напрямів прикладного мовознавства чим далі, тим помітніше місце займає корпусна лінгвістика. Увага до укладання й використання лінгвістичних корпусів зумовлена, з одного боку, тим, що корпус текстів – це потужна матеріальна й інструментальна база для різноманітних наукових та практичних робіт, а з іншого, – розвитком інформаційних технологій, які суттєво спрощують процедуру створення великих колекцій лінгвістичних даних.

Мовні корпуси дозволяють вирішувати різноаспектні лінгвістичні завдання: вивчати й описувати лексичні значення на базі реального мовленнєвого слововжитку [1, с. 44]; аналізувати граматичну структуру мови, верифікуючи теоретичні знання про неї [3, с. 224]; проводити дослідження з мовної стилістики [2, с. 90; 4]; прагматики [13, с. 5]; психо- та когнітивної лінгвістики [10, с. 23]; ілюструвати навчальні курси лінгвістичних дисциплін тощо. Особливим різновидом корпусних проектів є паралельні корпуси, що містять оригінали й переклади текстів певними мовами. З їх допомогою можна, як зазначає М. Шведова, “швидко отримувати велику кількість реальних перекладацьких рішень, що були прийняті носіями мови при створенні перекладу, й аналізувати виявлені відповідники в лексиці та граматиці, досліджувати перекладні моделі” [9, с. 100–102]. Паралельний корпус надає широкі можливості для статистично верифікованого дослідження міжмовних лексичних відповідностей з метою уточнення значень слів або кореляцій між їх уживанням у певних значеннях [7, с. 81]. Крім того, паралельні корпуси є безцінним матеріалом для навчання систем машинного перекладу [14, с. 41], а також можуть використовуватись як практичні довідники для самостійного вивчення іноземних мов і написання іншомовних текстів [11, с. 219].

Отже, потреба в паралельних корпусах цілком очевидна як для науковців, так і для широкого загалу користувачів.

Українська корпусна лінгвістика розвивається вже не перше десятиріччя і має вагомий досягнення. Сьогодні існують такі корпусні продукти: Український національний лінгвістичний корпус (Режим доступу: http://unlc.icybcluster.org.ua/virt_unlc/, розробник – Український мовно-інформаційний фонд НАН України), Корпус текстів української мови (Режим доступу: <http://corpora.pp.ua>, розробник кафедра української мови і прикладної лінгвістики Донецького національного університету), Корпус української мови (Режим доступу: <http://mova.info/corpus.aspx>, розробник лабораторія комп'ютерної лінгвістики Київського національного університету імені Тараса Шевченка) та ін. Проте, попри очевидний прогрес у цій галузі, досі лишається майже неопрацьованим завдання розробки паралельних українсько-іноземних корпусів текстів. Серед наявних нині ресурсів можна згадати паралельний російсько-український підкорпус Національного корпусу російської мови (Режим доступу: <http://ruscorpora.ru/search-para-uk.html>) [8], польсько-український корпус (Режим доступу: <http://domeczek.pl/~polukr>) [12], багатомовні корпуси ParaSol (Режим доступу: <http://parasolcorpus.org/ParaVoz/>) [15] та Intercorp (Режим доступу: https://kontext.korpus.cz/first_form?corpname=intercorp_v9_uk). Також ведеться робота над болгарсько-українським паралельним корпусом [7]. Відкритий доступ на сьогодні є тільки до російського та польського ресурсів. Оскільки паралельний корпус – необхідний інструмент для сучасного лінгвіста, перекладача, редактора, то в лабораторії комп'ютерної лінгвістики Інституту філології Київського національного університету імені Тараса Шевченка було започатковано проект зі створення багатомовного паралельного корпусу на базі Корпусу української мови (далі – ПарКУМ).

Мета статті – спираючись на досвід розробки проекту ПарКУМ, сформулювати основні засади створення та використання паралельного корпусу.

Завдання:

- визначити напрями перекладу, що складатимуть базу паралельного корпусу, а також окреслити принципи добору текстів;
- вибрати основні параметри розмітки;
- визначити концепції наповнення корпусу (статична vs динамічна структура, відкритість чи закритість системи редагування тощо);
- розробити структуру інтерфейсу користувача.

Одним із важливих питань, яке потрібно було вирішити розробникам ПарКУМ, став добір мовного матеріалу. Перший етап наповнення корпусу орієнтований на дві іноземні мови: англійську та російську; надалі планується доповнити базу паралельними текстами з турецької, польської, болгарської, німецької, білоруської та інших мов. Оскільки кожен напрям перекладу має свої специфіку та аспекти дослідження, вирішено було вводити до корпусу як українські версії іншомовних текстів, так і переклади україномовних матеріалів іншими мовами. Переклади українською або англійською мовою за посередництва інших мов, а також випадки паралельного перекладу українською та англійською іншомовних текстів поки що було вирішено не залучати.

Поетичні переклади до паралельного корпусу не включаються через специфіку художнього перекладу поезії (вільне поводження з лексикою, часом істотний відхід від змісту, а то й від форми оригіналу). Вони увійдуть (як самостійні одиниці) до поетичного підкорпусу КУМ.

У ПарКУМ реалізовано три види розмітки: метатекстову, структурну і власне лінгвістичну.

Метатекстова розмітка передбачає опис завантажених текстів за низкою параметрів. Вона включає інформацію про час створення тексту (точна дата або приблизний діапазон, якщо точна дата невідома); бібліографічний опис видання тексту (якщо це відомо) – ці дані вносяться як стосовно оригіналу, так і для перекладного тексту. Так само зазначаються дані про особу автора та перекладача (ім'я та прізвище, рік народження). Надалі планується доповнити метатекстовий опис жанрово-стильовими характеристиками. Метатекстова розмітка є важливим блоком корпусної розмітки, що виконує низку функцій, зокрема сприяє вибудові архітектури корпусу; дозволяє контролювати його наповнення, стежити за збалансованістю складу; надає користувачеві можливість добирати та групувати текстовий матеріал за різноманітними параметрами, закодованими розміткою [5, с. 62]: за хронологією, мовою оригіналу чи перекладу, за ім'ям автора чи навіть за його статтю. Це допоможе здійснювати пошук мовних явищ у сфері перекладу, обмежених цими параметрами (приміром, простежити, як передавали певне слово чи словосполучення в середині ХХ ст. та на початку ХХІ ст.). Таким чином, метатекстова розмітка допомагає дослідникові врахувати під час роботи з матеріалом не лише власне лінгвістичні, але й позамовні аспекти (історико-культурні, гендерні, функціонально-стилістичні тощо). Такий комплексний підхід до аналізу перекладної практики є важливим чинником у формуванні професійної культури перекладу текстів.

Структурна розмітка відіграє в паралельному корпусі особливо важливу роль, оскільки забезпечує зіставлення фрагментів оригіналу з відповідними фрагментами перекладного тексту. Вона не лише відкриває можливості до пошуку потрібних текстових елементів, а й унаочнює відмінності в лінгвістичній структурі й стилістиці оригінального та перекладного текстів. Структурна розмітка в корпусі здійснюється в два етапи. На першому етапі програма автоматично опрацьовує завантажені тексти, розбиваючи їх на речення. Зрозуміло, що пунктуаційні межі зіставляваних фрагментів не завжди тотожні, особливо у художніх текстах (одному реченню в оригіналі можуть відповідати два чи більше речень перекладу, і навпаки). Тому цілком імовірними є помилки, спричинені як розбіжністю структури текстів, так і характером роботи програми. З огляду на це, у структурному анотуванні передбачено ще один етап, коли користувач в автоматизованому режимі може виправляти помилки вирівнювання текстів (Рис. 1).

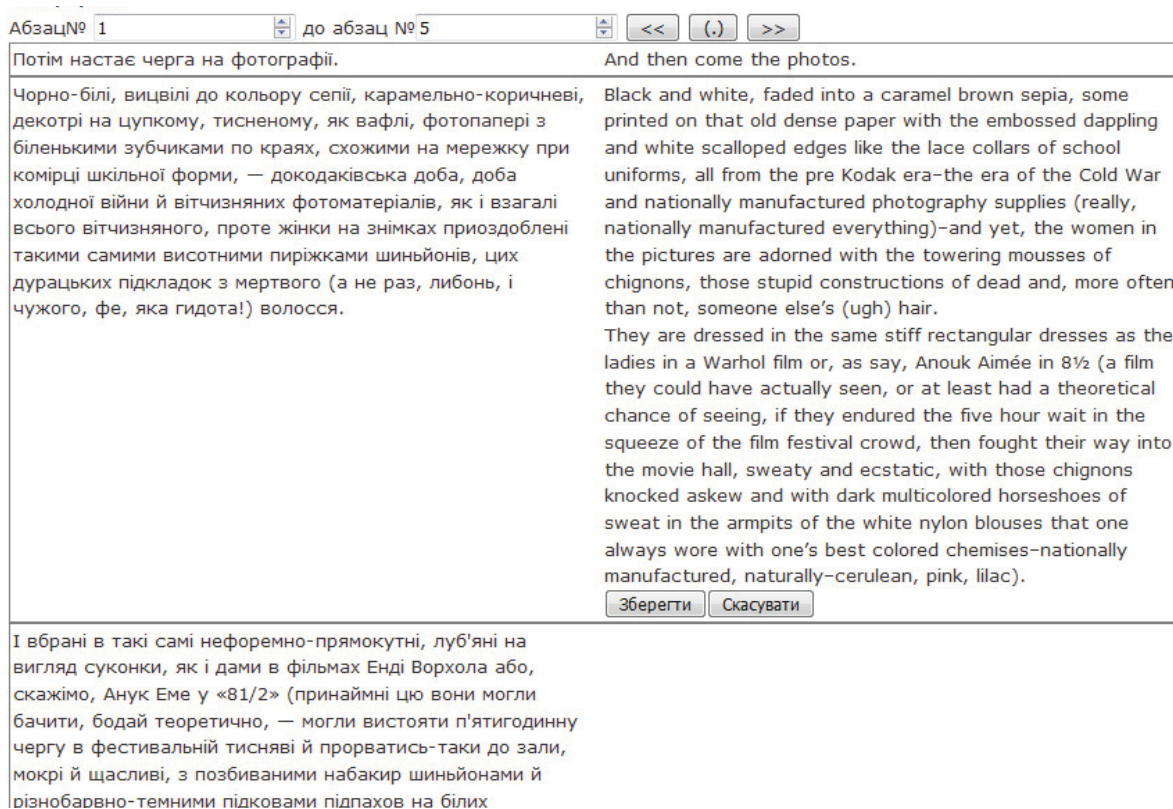


Рис. 1. Вирівнювання абзаців та речень

Загальний алгоритм уведення текстів у систему такий. Натискаючи на сторінці переліку текстів кнопку “Додати”, користувач переходить до форми метаопису, заповнивши яку, отримує можливість завантажити текст; для цього необхідно зазначити його мову, вказати шлях до файлу та перевірити коректність завантаження (див. Рис. 2); кожен із паралельних текстів завантажується окремим файлом.

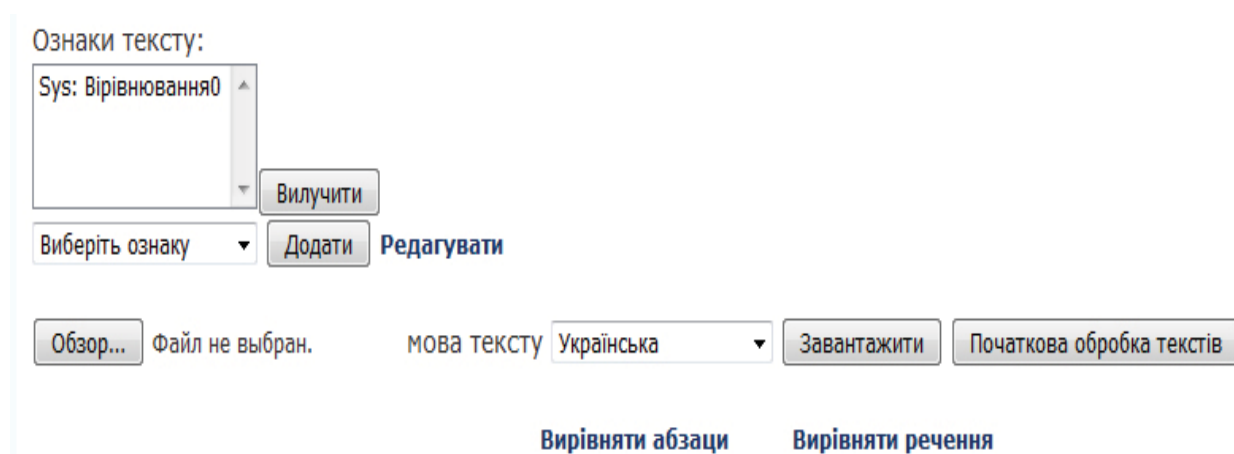


Рис. 2. Завантаження та обробка тексту

Наступним кроком є вирівнювання паралельних текстів, яке здійснюється на двох структурних рівнях – абзацу й речення; вирівнювання речень відбувається після того, як усунуто розбіжність на рівні абзаців.

Після уведення та обробки тексти передаються до лінгвістичного аналізатора для здійснення лінгвістичної розмітки.

На сьогодні в корпусі реалізовано морфологічну розмітку українського тексту, що дозволяє в разі обрання українсько-іноземного напрямку здійснювати пошук не лише за конкретними словами, а й за їх граматичними ознаками. Ця опція дозволяє користувачам порівнювати морфологічну структуру мов, добираючи ілюстрації їхніх спільних та відмінних рис чи розглядаючи особливості перекладу специфічних мовних явищ. Наступним кроком буде здійснення морфологічної розмітки для іншомовних (передусім, англійських) текстів.

Крім того, в корпусі заплановано уведення синтаксичної анотації для обох напрямів перекладу. Якщо структурна розмітка дозволяє порівняти обсяг текстових відповідників оригіналу й перекладу, то представлення синтаксичної структури фрагментів уможливить аналіз синтаксичних трансформацій, що виникли під час перекладу. Синтаксис українських текстів буде представлено двома підрівнями за зразком синтаксичного підкорпусу Корпусу української мови (Режим доступу: http://www.mova.info/syntaxis_search.aspx). Для словосполучень надаватиметься інформація про частиномовний тип сполуки та характер зв'язку в ньому. Структура речень подаватиметься у вигляді дерев залежностей – формалізованих схем речень, оформлених у вигляді математичних графів з вершинним членом-присудком [6, с. 8–9]. Накопичення синтаксичних моделей сприятиме формуванню бази типових перекладних трансформацій із широким ілюстративним матеріалом. Це уможливить зіставлення граматичної стилістики паралельних текстів. Застосування методу граматики залежностей для порівняльного опису синтаксичної структури оригінального й перекладного текстів дозволяє простежити, наскільки перекладач зміг передати особливості авторського синтаксису.

У роботі з наповнення ПарКУМ було вирішено послуговуватися принципом відкритого доступу до адміністрування корпусу. Проект розроблено як динамічну систему, в якій користувачам надається можливість самостійно поповнювати базу текстів та редагувати структуру й параметри уведеного матеріалу. Доступ до сторінки завантаження текстів і редагування метатекстової інформації надається через реєстрацію на порталі Mova.info. Це, з одного боку, пришвидшує процес наповнення корпусу, а з іншого – персоналізує пошук, оскільки користувач отримує можливість завантажити в систему потрібні йому тексти.

Інтерфейс пошуку в корпусі також орієнтований на максимальну персоналізацію (Рис. 3). Добір текстів можна здійснювати, зазначивши мовну пару й напрям перекладу та вказавши як одиницю пошуку слово або граматичну ознаку. Також доступна функція вибору конкретного тексту; тобто користувач може створити власний підкорпус, обравши тексти зі списку наявних або додавши до корпусу нові.

Вихід

Google Корист Пошук Вибрати м

TRANSLITERAZIYA SLOVNIKI PROEKTI CHITALNYA ZALA POSILANNYA

Напрямок: Ukr-Eng

Пошук за: паралельним текстом

Слово: it

Обмежити пошук в: всі

Пошук за морфологічною ознакою

Частина мови: всі

показати

| Невеличка драма | A Little Touch of Drama |
|--|--|
| В чарах кохання моє дівування, вільною пташкою хотіла б прожити, — вільне кохання і вільне обрання, серденьку воля, як хоче любить... | To spend my girlhood deep in love, I'd like to live like a free bird, all choices open to my heart which knows best what it craves. |

Рис. 3. Інтерфейс пошуку

Таким чином, ПарКУМ спроектовано як багатофункціональну довідкову систему з різними рівнями лінгвістичної анотації та динамічним наповненням. На сьогодні корпус має обмежену функціональність і працює в тестовому режимі. Ведеться активна робота з наповнення бази текстів українсько-англійської пари. Наступними кроками заплановано розширення функціональних можливостей за рахунок уведення додаткової лінгвістичної анотації, а також залучення до корпусу паралельних текстів іншими мовами.

Література:

1. Бобкова Т. В. Концепція колокації: корпусний підхід / Т. В. Бобкова // Науковий вісник Міжнародного гуманітарного університету. Серія: Філологія. – 2014. – Вип. 10 (2). – С. 42–45.
2. Бук С. Статистичні характеристики роману Івана Франка “Основи суспільності” / С. Бук // Вісник Національного університету “Львівська політехніка”. – Львів, 2010. – № 676. – С. 90–93.
3. Данилюк І. Г. Корпус текстів для вивчення граматичної службовості / І. Г. Данилюк // Лінгвістичні студії. – 2013. – Вип. 26. – С. 224–230.
4. Лангенбах М. О. Автоматизація стилістичних досліджень українських текстів / М. О. Лангенбах // Сучасна україністика: проблеми мови, літератури і культури: Оломоуцький симпозиум україністів. – Оломоуц, 2016. – Вип. 7. – С. 146–152.
5. Савчук С. О. Метатекстовая разметка в национальном корпусе русского языка: базовые принципы и основные функции [Электронный ресурс] / С. О. Савчук // Национальный корпус русского языка: 2003–2005. Результаты и перспективы. – Режим доступа : <http://ruscorp.org.ru/sbornik2005/05savchuk.pdf>
6. Севбо И. П. Графическое представление синтаксических структур и стилистическая диагностика / И. П. Севбо. – К. : Наукова думка, 1981. – 192 с.
7. Сірук О. Б. Лексичні перекладні еквіваленти в болгарських і українських паралельних текстах / О. Б. Сірук // Українське мовознавство. – К., 2013. – Вип. 43. – С. 75–86.

8. *Тищенко-Монастирська О.* Паралельні українсько-російський та російсько-український корпуси / *О. Тищенко-Монастирська, М. Шведова, Д. Січинава* // Лексикографічний бюлетень : [зб. наук. пр.]. – К. : Ін-т української мови НАН України, 2011. – Вип. 20. – С. 35–38.
9. *Шведова М.* Корпусна лінгвістика та лексико-граматична типологія / *М. Шведова, Д. Січинава* // Українське мовознавство. – К., 2013. – Вип. 43. – С. 95–103.
10. *Church K. W.* Word association norms, mutual information, and lexicography / *K. W. Church, N. J. Hanks* // *Computational Linguistics*, 2010. – No. 16. – P. 22–29.
11. *Frankenberg-Garcia A.* Lost in parallel concordances / *A. Frankenberg-Garcia* // *Corpora and Language Learners*. – Amsterdam-Philadelphia, 1996. – P. 213–232.
12. *Kotsyba N.* Polsko-Ukraiński Korpus Równoległy PolUKR i jego następcą PolUKR-2 / *N. Kotsyba* // *Polskojęzyczne korpusy równoległe*. – Warszawa, 2016. – P. 133–142.
13. *Pragmatics and Corpus Linguistics: A Mutualistic Entente* / [ed. *J. Romero-Trillo*]. – Berlin : Mouton de Gruyter, 2008 – 282 p.
14. *Utiyama M.* Mining patterns from parallel corpora / *M. Utiyama, H. Isahara* // *Learning Machine Translation*. – Cambridge, London : The MIT Press, 2009. – P. 41–58.
15. *Waldenfels R.* ParaSol: Introduction to a Slavic Parallel Corpus / *R. von Waldenfels* // *Prace Filologiczne*. – Warszawa : Wydział Polonistyki Uniwersytetu Warszawskiego, 2012. – No. LXIII. – P. 293–302.

References:

1. *Bobkova T. V.* Kontsepsiia kolokatsii: korpusnyi pidkhid [The conception of the collocation: the corpus approach] / *T. V. Bobkova* // *Naukovyi visnyk Mizhnarodnoho humanitarnoho universytetu. Serii: Filolohiia*. – 2014. – Vyp. 10 (2). – S. 42–45.
2. *Buk S.* Statystychni kharakterystyky romanu Ivana Franka “Osnovy suspilnosti” [The statistical characteristics of the Ivan Franko’s novel “The Basis of Publicness”] / *S. Buk* // *Visnyk Natsionalnoho universytetu “Lvivska politekhnika”*. – Lviv, 2010. – No. 676. – S. 90–93.
3. *Danyliuk I. H.* Korpus tekstiv dlia vyvchennia hramatychnoi sluzhbovosti [The textual corpus for the grammatical functionality study] / *I. H. Danyliuk* // *Linhvistychni studii*. – 2013. – Vyp. 26. – S. 224–230.
4. *Langenbakh M. O.* Avtomatyzatsiia stylistychnykh doslidzhen ukrainskykh tekstiv [The automatization of the stylistic studies of the Ukrainian texts] / *M. O. Langenbakh* // *Suchasna ukrainistyka: problemy movy, literatury i kultury: Olomoutskyi sympozium ukrainistiv*. – Olomouts, 2016. – Vyp. 7. – S. 146–152.
5. *Savchuk S. O.* Metatekstovaya razmetka v natsionalnom korpuse russkogo yazyka: bazovye printsipy i osnovnye funktsii [The metadata in the annotation of The National Corpus of the Russian Language] / *S. O. Savchuk* // *Natsionalnyy korpus russkogo yazyka: 2003–2005. Rezultaty i perspektivy*. – Rezhim dostupu : <http://ruscorpora.ru/sbornik2005/05savchuk.pdf>
6. *Siruk O. B.* Leksychni perekladni ekvivalenty v bolharskykh i ukrainskykh paralelnykh tekstakh [The lexical translation equivalents in the Bulgarian and Ukrainian texts] / *O. B. Siruk* // *Ukrainske movoznavstvo*. – K., 2013. – No. 43. – S. 75–86.
7. *Sevbo I. P.* Graficheskoe predstavlenie sintaksicheskikh struktur i stilisticheskaya diagnostika [The Graphic representation of the syntactic structures and the stylistic diagnostics] / *I. P. Sevbo*. – K. : Naukova dumka, 1981. – 192 s.
8. *Tyshchenko-Monastyrska O.* Paralelni ukrainsko-rosiiskyi ta rosiisko-ukrainskyi korpusy [The parallel Russian-Ukrainian and Ukrainian-Russian corpora] / *O. Tyshchenko-Monastyrska, M. Shvedova, D. Sichinava* // *Leksykohrafichnyi biuleten : [zb. nauk. pr.]*. – K. : Ін-т української мови НАН України, 2011. – Вип. 20. – С. 35–38.
9. *Shvedova M.* Korpusna linhvistyka ta leksyko-hramatychna typolohiia [The corpus linguistics and the lexical and grammatical typology] / *M. Shvedova, D. Sichinava* // *Ukrainske movoznavstvo*. – K., 2013. – Vyp. 43. – S. 95–103.
10. *Church K. W.* Word association norms, mutual information, and lexicography / *K. W. Church, N. J. Hanks* // *Computational Linguistics*, 2010. – No. 16. – P. 22–29.
11. *Frankenberg-Garcia A.* Lost in parallel concordances / *A. Frankenberg-Garcia* // *Corpora and Language Learners*. – Amsterdam-Philadelphia, 1996. – P. 213–232.
12. *Kotsyba N.* Polsko-Ukraiński Korpus Równoległy PolUKR i jego następcą PolUKR-2 / *N. Kotsyba* // *Polskojęzyczne korpusy równoległe*. – Warszawa, 2016. – P. 133–142.
13. *Pragmatics and Corpus Linguistics: A Mutualistic Entente* / [ed. *J. Romero-Trillo*]. – Berlin : Mouton de Gruyter, 2008 – 282 p.

14. *Utiyama M. Mining patterns from parallel corpora / M. Utiyama, H. Isahara // Learning Machine Translation. – Cambridge, London : The MIT Press, 2009. – P. 41–58.*
15. *Waldenfels R. ParaSol: Introduction to a Slavic Parallel Corpus / R. von Waldenfels // Prace Filologiczne. – Warszawa : Wydział Polonistyki Uniwersytetu Warszawskiego, 2012. – No. LXIII. – P. 293–302.*

Дарчук Н. П., Лангенбах М. О., Сорокин В. М., Ходаковская Я. В. Параллельный корпус текстов ПарКУМ.

Несмотря на активное развитие корпусной лингвистики в Украине, до сих пор существует большой пробел в области разработки параллельных корпусов. На сегодня почти не существует в открытом доступе подобных проектов, содержащих украиноязычные тексты. Цель статьи – изложение основных принципов создания и использования корпуса параллельных текстов ПарКУМ. Задания, решаемые в ходе исследования: определение принципов подбора текстов; выбор основных параметров разметки; формирование концепции работы с материалом; разработка структуры пользовательского интерфейса. Предоставляется информация обо всех типах разметки, предусмотренных в корпусе: метатекстовой, структурной и лингвистической, обеспечивающей возможность поиска информации не только по конкретным лексемам, но и по грамматическим признакам.

Ключевые слова: параллельный корпус, корпусная лингвистика, аннотация корпуса, параллельные тексты, переводные соответствия.

Darchuk N. P., Langenbakh M. O., Sorokin V. M., Khodakivska Ya. V. Parallel corpus of texts ПарКУМ.

Despite the fact that Ukrainian corpus linguistics has some visible achievements, there is one field which is still almost unexplored – parallel corpora. In this paper we present a project of parallel corpus containing Ukrainian texts. The goal of our research is to formulate the basic principles of parallel corpus development. The tasks being solved are: to define the directions of translation and demands to the textual material; to choose necessary parameters of annotation; to build the architecture of corpus system and define user roles; to develop user interface. The article gives the information about all types of tagging specified for the corpus texts: metadata, structural and linguistic annotation. The corpus works in two modes: administrative (available after the registration at the <http://mova.info>) and search. The project works in test mode. The Ukrainian-English and English-Ukrainian parallel texts are being collected now and some examples of them are already available for search. On the next stages the corpus will be filled with other parallel texts – Polish, Bulgarian, Turkish, German etc.

Keywords: parallel corpus, corpus linguistics, annotation of the corpus, parallel texts, translation equivalents.

УДК 81'37:[811.161.2+811.162.1+811.111

Деменчук О. В.

Рівненський державний гуманітарний університет

МОДЕЛІ СЕМАНТИЧНОЇ ДЕРИВАЦІЇ ІРРАЦІОНАЛЬНОЇ ЛЕКСИКИ

У розвідці схарактеризовано моделі семантичної деривації ірраціональної лексики – семантичного класу слів, які позначають стан людини, що виникає не на раціональній основі, ґрунтується не на вимогах розуму, логіки, суджень і т. ін. Основну увагу приділено визначенню моделей ірраціональної ад'єктивної лексики в українській, польській та англійській мовах, розкриттю зв'язків між вихідним та цільовим значеннями, з'ясуванню динаміки розвитку семантичної парадигми названого класу лексики в зіставному аспекті. Установлено, що у процесі семантичної деривації ірраціональна ад'єктивна лексика виявляє ознаки похідності ситуативного типу, пов'язані зі змінами характеристик учасників ситуації.

Ключові слова: ірраціональна лексика, модель, семантична деривація, учасник, концепт ситуації, дериваційні відношення.