

*Валентина Балог, к. ф. н.\**

*Орися Демська-Кульчицька, к. ф. н.\**

Інститут української мови НАН України (Київ)

УДК 161.2.81'374.72'22

## **НАЦІОНАЛЬНИЙ КОРПУС УКРАЇНСЬКОЇ МОВИ: ПІДКОРПУС ХУДОЖНЬОЇ ЛІТЕРАТУРИ**

*У статті висвітлено засади формування та формалізації матеріалів підкорпусу художньої літератури в складі генерального корпусу НКУМ: обґрунтовано принципи добору текстів, визначено кількість та обсяг текстових фрагментів, подано адаптовану систему кодування первинних даних.*

Створення *Національного корпусу української мови* [1] є першою спробою застосування загальної теорії корпусної лінгвістики до української мови, наслідком чого має стати стандартний загальномовний корпус, зіставний з аналогічними світовими загальномовними текстовими корпусами національного типу. Такий підхід, на наш погляд, повинен розв'язати низку завдань, найважливішими з яких є впровадження корпусно базованих методик дослідження української мови, формування корпусної лінгвоукраїністики як окремого напрямку, а також забезпечення технологічного статусу української мови в умовах інформаційного суспільства.

Національний корпус – це зібрання текстів, що репрезентують національну мову на певному(их) етапі(ях) її існування в усьому різноманітті жанрів, стилів, історичних, територіальних і соціальних варіантів. Більшість мов світу й усі європейські мови на сьогодні вже мають власні національні корпуси, які статусом прирівнюють до

---

\* © В.Балог, 2005

\* © О. Демська-Кульчицька, 2005

академічних граматик та великих тлумачних словників національних мов. Відповідно до цієї дефініції сформулюємо визначення поняття 'національний корпус української мови'.

**Національний корпус української мови (НКУМ)** – це перетворена на електронну форму, стандартно організована й програмно оброблена вибірка текстів сучасної української мови, репрезентативна для всіх функціональних рівнів загальнонародної мови, призначена для лінгвістичного аналізу й технологічного застосування.

Стосовно лінгвістичного аналізу, то йдеться про академічні дослідження різних рівнів мовної системи (лексики, фразеології, термінології, морфології та синтаксису, словотвору, орфографії тощо), про методику української мови як рідної й як іноземної тощо. Натомість, технологічне застосування передбачає використання корпусу з метою побудови машинної мовної моделі як основи для розробок у галузі інформаційних технологій, створення програм автоматичного розпізнавання й синтезу українського мовлення, забезпечення автоматичних методів перетворення українськомовної текстової інформації, лінгвістичної підтримки автоматичних систем управління.

З погляду типологічно-аплікативних характеристик НКУМ розглядаємо як

- **дослідницький**: орієнтований на широкий клас передусім лінгвістичних завдань;
- **фрагментний**: збудований з текстових фрагментів чи уривків текстів, відібраних за попередньо визначеними засадами відбору текстових даних до корпусу;
- **мішаний**: передбачає введення текстових фрагментів обох варіантів реалізації мовної системи – писемних і усних;
- **динамічний**: передбачає константне поповнення множини корпусних текстів;
- **синхронний**: охоплює рівень сучасної української мови;
- **загальнонародної (національної) мови**: з урахуванням територіальної специфіки як у межах України, так і за її межами;
- **мономовний**: тексти, що увійшли до корпусу, є результатом мовної діяльності носіїв української мови;

- **морфологічно анотований**: усі текстові дані розмічені до рівня слова й кожне слово передбачає маркування частиномовної належності та відповідних морфологічних значень.

Побудова довільного корпусу на етапі проекту вимагає обов'язкової параметризації предметної галузі, репрезентованої корпусом. Оскільки національний корпус є зібранням електронних текстів, завдання яких репрезентувати національну мову на певному етапі(ах) її існування в усьому різноманітті жанрів, стилів, історичних, територіальних і соціальних варіантів, то предметну галузь НКУМ повинна становити українська загальнонародна мова в таких формах її існування, як літературна мова, територіальний діалект і соціолект. Найповніше в НКУМ повинна бути репрезентована літературна мова – основна наддіалектна форма існування природної мови, ознаками якої є опрацьованість, унормованість, поліфункціональність, стилістична диференціація, фіксованість.

Змістова параметризація предметної галузі НКУМ передбачає визначення джерельної текстової бази корпусу з наступною стилістично-жанровою стратифікацією відібраних текстів. Тут стильова диференціація мови можлива, по-перше, за характером мовної експресії, і, по-друге, за характером суспільної функції мови.

За характером мовної експресії традиційно розрізняють високий, середній і низький стиль. В основі такої тристильової диференціації лежить концепція залежності між предметом викладу, тематикою і добором мовних засобів та жанрів. Цей поділ, успадкований європейською традицією доби Відродження й бароко з александрійської філософської школи античного періоду [2, с. 39], детермінував жанровий підхід до параметризації предметного домену низки національних корпусів європейських мов.

В українській традиції тристильова диференціація мови, яку, зокрема, розробляли Ф. Прокопович, М. Довгалевський, Г. Кониський, втратила актуальність у процесі історичного розвитку стилістичної системи української мови, коли три традиційні стилі ('слоги') занепали, а основними одиницями стильової диференціації української мови стали структурно-функціональні стилі. Таким чином, в українську лінгвістичну традицію увійшов функціональний принцип стилістичної

диференціації мови, і ця диференціація покладена в основу змістової параметризації НКУМ.

Приймаємо, що сучасна українська мова на рівні стильової диференціації вкладається в семиелементну систему: 1) художній, 2) науковий, 3) офіційно-діловий, 4) публіцистичний, 5) конфесійний, 6) розмовний і 7) епістолярний стилі. А залежно від реалізації мовної системи виділяємо тексти усного й писемного варіантів національної мови. Художній, науковий, офіційно-діловий і публіцистичний стилі в межах НКУМ головно репрезентовані писемними текстами, теологічний – писемними і усними, а розмовний – усними.

Значимо, що корпус текстів лише тоді відповідатиме загальноприйнятим вимогам корпусу національного типу і виконуватиме всі передбачені для нього функції, коли в ньому адекватно репрезентовано обидві форми функціонування мовної системи. Але для багатьох галузей мовної практики писемна форма національної мови є пріоритетною і вона зазвичай реалізує себе в художньому, науковому, офіційно-діловому і публіцистичному стилях. У НКУМ приймаємо, що фактичний матеріал писемного варіанту української мови становитимуть тексти таких функціональних стилів:

- 1) художнього;
- 2) наукового;
- 3) офіційно-ділового;
- 4) публіцистичного;
- 5) конфесійного;
- 6) епістолярного.

Художній стиль у НКУМ репрезентуватимуть прозові, поетичні й драматургічні тексти відповідних хронологічних меж (підкорпуси сучасної української мови й української мови зламу ХХ / ХХІ ст.) з такою жанровою диференціацією:

- проза: роман, повість, оповідання, новела;
- поезія: вірш, сонет, поема, балада, сатира;
- драматургія: драма, комедія, трагедія.

Загальний обсяг корпусу детермінує обсяг його складових як засіб досягнення збалансованості корпусу. Незбалансованість дослідницького

корпусу національного типу призводить до неможливості проведення коректних досліджень на його основі, оскільки такий корпус не відтворює загальної мовної картини через відхилення в бік того чи іншого мовного матеріалу. Наприклад, якщо не збалансувати дані усної і писемної форми мови, то за умови перенасичення даних усними текстами матимемо неадекватну картину функціонування норм. Отже, пропонуємо такі загальнокорпусні кількісні параметри НКУМ (див. табл. 1).

Таблиця 1. Загальні кількісні параметри НКУМ

характер текстових даних	відсоток від підкорпусу	кількість слововживань
<b>ПИСЕМНІ ТЕКСТИ</b> (850 тис. слів – 85 %)		
<i><b>Художній стиль (350 тис. слів – 35%)</b></i>		
<i>Проза (250 тис. слів – 25%)</i>		
роман	10%	100 тис.
повість	10%	100 тис.
оповідання	3%	30 тис.
новела	2%	20 тис.
<i>Поезія (50 тис. слів – 5%)</i>		
вірш (лірика)	1%	10 тис.
сонет	1%	10 тис.
поема	1%	10 тис.
балада	1%	10 тис.
сатира	1%	10 тис.
<i>Драматургія (50 тис. слів – 5%)</i>		
драма	2%	20 тис.
комедія	1,6%	16 тис.
трагедія	1,4%	14 тис.
<i><b>Науковий стиль (100 тис. слів – 10%)</b></i>		
<i>Власне науковий підстиль (20 тис. слів – 2%)</i>		
тексти гуманітарних наук	1%	10 тис.

тексти природничих і точних наук	1%	10 тис.
<i>Науково-популярний підстиль (20 тис. слів – 2%)</i>		
тексти гуманітарних наук	1%	10 тис.
тексти природничих і точних наук	1%	10 тис.
<i>Науково-методичний підстиль (20 тис. слів – 2%)</i>		
тексти гуманітарних наук	1%	10 тис.
тексти природничих і точних наук	1%	10 тис.
<i>Професійно-технічний підстиль (40 тис. слів – 4%)</i>		
культура і мистецтво	1%	10 тис.
міжнародні відносини	0,5%	5 тис.
бізнес	1%	10 тис.
медицина	1%	10 тис.
техніка	0,5%	5 тис.
<i>Офіційно-діловий стиль (50 тис. слів – 5%)</i>		
законодавчо-правові державні документи	2,5%	25 тис.
організаційно-службові документи	1%	10 тис.
суспільна документація	1,5%	15 тис.
<i>Публіцистичний стиль (250 тис. слів – 25%)</i>		
періодика	24%	240 тис.
візуальна реклама	0,5%	5 тис.
учнівські та студентські неопубліковані твори	0,5%	5 тис.
<i>Конфесійний стиль (50 тис. слів – 5%)</i>		
Біблія	3%	30 тис.

твори літургійного призначення	2%	20 тис.
<b>Епістолярний стиль (50 тис. слів – 5%)</b>		
<i>Класична кореспонденція (40 тис. слів – 4%)</i>		
офіційно-ділове внутрішньодержавне листування	1%	10 тис.
офіційне міждержавне листування	1%	10 тис.
приватне листування	2%	20 тис.
<i>Електронна кореспонденція (10 тис. слів – 1%)</i>		
е-mail листування	1%	10 тис.
УСНІ ТЕКСТИ (150 тис. слів – 15 %)		
<b>Офіційно-діловий стиль (70 тис. слів – 7%)</b>		
живі службово- організаційні розмови	2%	20 тис.
телефонні службово- організаційні розмови	1%	10 тис.
виступи, повідомлення, оголошення на зборах	4%	40 тис.
<b>Розмовний стиль (50 тис. слів – 5%)</b>		
побутові розповіді	1.5%	15 тис.
побутове діалогічне мовлення	0,5%	5 тис.
повідомлення в транспорті	0,5%	5 тис.
аудіореклама	1%	10 тис.
телефонна побутова комунікація	0,5%	5 тис.
діалектне мовлення	1%	10 тис.

<i>Конфесійний стиль (30 тис. слів – 3%)</i>		
проповіді	3%	30 тис.
Всього підкорпус:	100%	1 млн.

Загальний обсяг НКУМ становитиме **1 млн.** слововживань. Очевидно, що згідно із запропонованою стратегією розвитку корпусу, ця величина є початковою чи вихідною.

Наступним кроком проектування статистики НКУМ є визначення кількісних параметрів текстових фрагментів кожного із функціональних стилів / підстилів і кількості слововживань у кожному конкретному текстовому фрагменті (див. табл. 2). Точкою відліку для цих розрахунків є загальна кількість слововживань генерального корпусу, у нашому випадку 1 млн. слововживань.

Таблиця 2. *Індивідуальні кількісні параметри тестових фрагментів НКУМ*

характер текстових даних	кількість фрагментів	кількість слів у фрагменті
Писемні тексти		
<i>Художній стиль</i>		
<i>Проза (250 тис. слів)</i>		
роман (100 тис. слів)	20	5 тис.
повість (100 тис. слів)	20	5 тис.
оповідання (30 тис. слів)	10	3 тис.
новела (20 тис. слів)	10	2 тис.
<i>Поезія (50 тис. слів)</i>		
вірш (лірика) (10 тис. слів)	20	500
сонет (10 тис. слів)	20	500
поема (10 тис. слів)	10	1 тис.
балада (10 тис. слів)	10	1 тис.
сатира (10 тис. слів)	20	500
<i>Драматургія (50 тис. слів)</i>		
драма (20 тис. слів)	10	2 тис.
комедія (16 тис. слів)	8	2 тис.
трагедія (14 тис. слів)	7	2 тис.



<b>Науковий стиль (100 тис. слів)</b>		
<i>Власне науковий підстиль (20 тис.)</i>		
тексти гуманітарних наук (10 тис. слів)	10	1 тис.
тексти природничих наук (10 тис. слів)	10	1 тис.
<i>Науково-популярний підстиль (20 тис. слів)</i>		
тексти гуманітарних наук (10 тис. слів)	10	1 тис.
тексти природничих наук (10 тис. слів)	10	1 тис.
<i>Науково-методичний підстиль (20 тис. слів)</i>		
тексти гуманітарних наук (10 тис. слів)	10	1 тис.
тексти природничих наук (10 тис. слів)	10	1 тис.
<i>Професійно-технічний підстиль (40 тис. слів)</i>		
культура і мистецтво (10 тис. слів)	10	1 тис.
міжнародні відносини (5 тис. слів)	5	1 тис.
бізнес (10 тис. слів)	10	1 тис.
медицина (10 тис. слів)	10	1 тис.
техніка (5 тис. слів)	5	1 тис.
<i>Офіційно-діловий стиль (50 тис. слів)</i>		
законодавчо-правові державні документи (25 тис. слів)	10	2.5 тис.
організаційно-службові документи (10 тис. слів)	20	500
суспільна документація (15 тис. слів)	10	1.5 тис.
<i>Публіцистичний стиль (250 тис. слів)</i>		
періодика (240 тис. слів)	120	2 тис.

візуальна реклама (5 тис. слів)	100	50
учнівські та студентські неопубліковані твори (5 тис. слів)	10	500
<b>Конфесійний стиль (50 тис. слів)</b>		
Біблія (30 тис. слів)	15	2 тис.
твори літургійного призначення (20 тис. слів)	10	2 тис.
<b>Епістолярний стиль (50 тис. слів)</b>		
<i>Класична кореспонденція (40 тис. слів)</i>		
офіційно-ділове внутрішньодержавне листування (10 тис. слів)	10	1 тис.
офіційне міждержавне листування (10 тис. слів)	10	1 тис.
приватне листування (20 тис. слів)	20	1 тис.
<i>Електронна кореспонденція (10 тис. слів)</i>		
е-mail листування (10 тис. слів)	20	500
Усні тексти		
<b>Офіційно-діловий стиль (70 тис. слів)</b>		
живі службово-організаційні розмови (20 тис. слів)	40	500
телефонні службово-організаційні розмови (10 тис. слів)	40	250
виступи, повідомлення, оголошення на зборах (40 тис. слів)	40	1 тис.

<i>Розмовний стиль (50 тис. слів)</i>		
побутові розповіді (15 тис. слів)	15	1 тис.
побутове діалогічне мовлення (5 тис. слів)	10	500
повідомлення в транспорті (5 тис. слів)	33	150
аудіореклама (10 тис. слів)	66	150
телефонна побутова комунікація (5 тис. слів)	20	250
діалектне мовлення (10 тис. слів)	10	1 тис.
<i>Конфесійний стиль (30 тис. слів)</i>		
проповіді (30 тис. слів)	15	2 тис.

Виходячи із загальних корпусних кількісних параметрів та індивідуальних кількісних параметрів тестових фрагментів, тексти художньої літератури мають становити 35% від усього корпусу, тобто якщо НКУМ плановано на 1 млн. слововживань, то обсяг художніх текстів становитиме 350 тис. слововживань, де, відповідно, проза – 250 тис. слововживань чи 35%, поезія – 50 тис. слововживань чи 5% і драматургія – 50 тис. слововживань чи 5%. Далі в межах прози, поезії та драматургії пропонуємо жанрову диференціацію текстів і, відповідно до визначених жанрів, забезпечення корпусу фактичним матеріалом: роман – 100 тис. слововживань чи 10%, повість – 100 тис. слововживань чи 10%, оповідання – 30 тис. слововживань чи 3%, новела – 20 тис. слововживань чи 2%; лірика – 10 тис. слововживань чи 1%, сонет – 10 тис. слововживань чи 1%, поема – 10 тис. слововживань чи 1%, балада – 10 тис. слововживань чи 1%, сатира – 10 тис. слововживань чи 1%; драма – 20 тис. слововживань чи 2%, комедія – 15 тис. слововживань чи 1,5% і трагедія – 15 тис. слововживань чи 1,5%.

Зважаючи на такі загальні корпусні кількісні параметри, індивідуальні кількісні параметри тестових фрагментів повинні становити для:

- ◆ прози – по 20 творів для роману і повісті, фрагменти яких повинні охоплювати по 5 тис. слововживань, і по 10 творів для оповідання й новели з фрагментами, відповідно, по 3 і 2 тис. слововживань;
- ◆ поезії – по 20 творів для лірики, сонету й сатири, з фрагментами по 500 слововживань, та по 10 творів для поеми і балади, з фрагментами по 1 тис. слововживань;
- ◆ драматургії – 10 творів з фрагментами на 2 тис. слововживань для драми, 8 творів на 2 тис. слововживань для комедії й 7 творів на 2 тис. слововживань для трагедії.

За аналогією до *Британського національного корпусу*, який розглядаємо як еталонний у сучасній корпусній лінгвістиці, вважаємо за доцільне, особливо на етапі впровадження напряму корпусної лінгвістики в українську мовознавчу традицію, визначити хронологічні межі НКУМ сучасною українською мовою. Тобто від останніх років XVIII ст. і до сьогодні, де окремо розглядати субперіод зламу XX / XXI ст., залишаючи поза увагою усі попередні періоди існування української мови, а саме: давньоукраїнський період від середини XI ст. до кінця XIV ст., ранньосередньоукраїнський від початку XV ст. до середини XVI ст., середньоукраїнський від середини XVI ст. до перших років XVIII ст., пізньосередньоукраїнський від середини і до кінця XVIII ст.

Виокремлення субперіоду зламу XX / XXI ст. мотивоване передусім екстралінгвальними чинниками, зокрема, утворенням Української держави і, відповідно, розширенням сфери функціонування української мови аж до входження в комп'ютерне середовище.

Отже, підкорпус художньої літератури у складі генерального корпусу НКУМ матиме обсяг 350 тис. слововживань, які репрезентуватимуть 155 текстів, і охоплюватиме художні твори сучасної української літературної мови від останніх років XVIII ст. і до сьогодні, де окремо слід розглядати субперіод зламу XX / XXI ст.

До укладання підкорпусу художньої літератури, як і до НКУМ загалом, пропонуємо застосувати принцип випадкової вибірки і

стосовно авторів, і стосовно творів. У результаті випадкової вибірки ми отримали такий список авторів, тексти яких становитимуть джерельну базу для підкорпусу художньої літератури:

- |  |   |
|--|---|
| 1. <a href="#">Андрій Головка</a>          | 35. <a href="#">Свген Маланюк</a>         |
| 2. <a href="#">Андрій Малишко</a>          | 36. <a href="#">Іван Багряний</a>         |
| 3. <a href="#">Андрій Чайковський</a>      | 37. <a href="#">Іван Гнатюк</a>           |
| 4. <a href="#">Архип Тесленко</a>          | 38. <a href="#">Іван Драч</a>             |
| 5. <a href="#">Богдан Лепкий</a>           | 39. <a href="#">Іван Карпенко-Карий</a>   |
| 6. <a href="#">Богдан-Ігор Антонич</a>     | 40. <a href="#">Іван Котляревський</a>    |
| 7. Борис Антоненко-Давидович               | 41. <a href="#">Іван Кочерга</a>          |
| 8. Борис Грінченко                         | 42. <a href="#">Іван Нечуй-Левицький</a>  |
| 9. Борис Олійник                           | 43. <a href="#">Іван Франко</a>           |
| 10. Валер'ян Підмогильний                  | 44. <a href="#">Катерина Мотрич</a>       |
| 11. Василь Барка                           | 45. <a href="#">Катря Гриневичева</a>     |
| 12. <a href="#">Василь Еллан Блакитний</a> | 46. Леонід Глібов                         |
| 13. <a href="#">Василь Земляк</a>          | 47. <a href="#">Леся Українка</a>         |
| 14. Василь Пачовський                      | 48. <a href="#">Ліна Костенко</a>         |
| 15. <a href="#">Василь Симоненко</a>       | 49. <a href="#">Максим Рильський</a>      |
| 16. <a href="#">Василь Стефаник</a>        | 50. <a href="#">Марійка Підгірянка</a>    |
| 17. <a href="#">Василь Стус</a>            | 51. <a href="#">Маркіян Шашкевич</a>      |
| 18. <a href="#">Віктор Неборак</a>         | 52. <a href="#">Марко Вовчок</a>          |
| 19. <a href="#">Віктор Петров</a>          | 53. <a href="#">Марко Кропивницький</a>   |
| 20. <a href="#">Володимир Винниченко</a>   | 54. <a href="#">Микола Бажан</a>          |
| 21. <a href="#">Володимир Діброва</a>      | 55. <a href="#">Микола Вінграновський</a> |
| 22. <a href="#">Володимир Дрозд</a>        | 56. <a href="#">Микола Вороний</a>        |
| 23. <a href="#">Володимир Малик</a>        | 57. <a href="#">Микола Зеров</a>          |
| 24. <a href="#">Володимир Самійленко</a>   | 58. <a href="#">Микола Куліш</a>          |
| 25. <a href="#">Володимир Сосюра</a>       | 59. <a href="#">Микола Олійник</a>        |
| 26. <a href="#">Всеволод Нестайко</a>      | 60. <a href="#">Микола Сингаївський</a>   |
| 27. В'ячеслав Сахно                        | 61. <a href="#">Микола Хвильовий</a>      |
| 28. <a href="#">Гнат Хоткевич</a>          | 62. Михайло Драй-Хмара                    |
| 29. Григорій Квітка-Основ'яненко           | 63. <a href="#">Михайло Коцюбинський</a>  |
| 30. <a href="#">Григорій Чубай</a>         | 64. <a href="#">Михайло Петренко</a>      |
| 31. <a href="#">Дмитро Білоус</a>          | 65. Михайло Старицький                    |
| 32. <a href="#">Дмитро Павличко</a>        | 66. <a href="#">Михайло Стельмах</a>      |
| 33. <a href="#">Емма Андіївська</a>        | 67. <a href="#">Михайль Семенко</a>       |
| 34. <a href="#">Свген Гуцало</a>           |   |

- 
- |  |  |
|--|--|
| 68. <a href="#">Наталена Королева</a>      | 89. <a href="#">Пантелеймон Куліш</a>        |
| 69. <a href="#">Наталка Білоцерківська</a> | 90. <a href="#">Петро Гулак-Артемівський</a> |
| 70. <a href="#">Оксана Забужко</a>         | 91. <a href="#">Платон Воронько</a>          |
| 71. <a href="#">Оксана Іваненко</a>        | 92. <a href="#">Роман Іваничук</a>           |
| 72. <a href="#">Олег Ольжич</a>            | 93. <a href="#">Семен Скляренко</a>          |
| 73. <a href="#">Олександр Довженко</a>     | 94. <a href="#">Сергій Плачинда</a>          |
| 74. <a href="#">Олександр Олесь</a>        | 95. <a href="#">Сидір Воробкевич</a>         |
| 75. <a href="#">Олександр Підсуха</a>      | 96. <a href="#">Степан Васильченко</a>       |
| 76. <a href="#">Олексій Коломієць</a>      | 97. <a href="#">Степан Руданський</a>        |
| 77. <a href="#">Олена Теліга</a>           | 98. <a href="#">Тарас Шевченко</a>           |
| 78. <a href="#">Олесь Бердник</a>          | 99. <a href="#">Тимофій Бордуляк</a>         |
| 79. <a href="#">Олесь Гончар</a>           | 100. <a href="#">Тодось Осьмачка</a>         |
| 80. <a href="#">Олеся Садова</a>           | 101. <a href="#">Улас Самчук</a>             |
| 81. <a href="#">Ольга Кобилянська</a>      | 102. <a href="#">Юрій Андрухович</a>         |
| 82. <a href="#">Осип Назарук</a>           | 103. <a href="#">Юрій Винничук</a>           |
| 83. <a href="#">Остап Вишня</a>            | 104. <a href="#">Юрій Клен</a>               |
| 84. <a href="#">Павло Глазовий</a>         | 105. <a href="#">Юрій Ліпа</a>               |
| 85. <a href="#">Павло Грабовський</a>      | 106. <a href="#">Юрій Мушкетик</a>           |
| 86. <a href="#">Павло Загребельний</a>     | 107. <a href="#">Юрій Покальчук</a>          |
| 87. <a href="#">Павло Тичина</a>           | 108. <a href="#">Юрій Федькович</a>          |
| 88. <a href="#">Панас Мирний</a>           | 109. <a href="#">Яр Славутич</a>             |

Оброблення текстів у НКУМ доцільно реалізовувати за допомогою засобів SGML [3] у форматі TEI [4]. Де SGML (Standard Generalized Markup Language) – це міжнародний стандарт на визначення незалежних від пристроїв уведення/виведення інформації, незалежних від обчислювального середовища методів подання текстів у електронній формі. А TEI (Text Encoding Initiative) – система кодування текстів, яка є міжнародним і міждисциплінарним стандартом подання усіх типів текстів, функціональних у бібліотечній, музейній, видавничій справах, лінгвістиці.

Схема кодування TEI використовує стандартну мову узагальненої розмітки (SGML) і спрямована на забезпечення обміну інформацією, що зберігається в електронній формі. Використання SGML і TEI – це шлях до універсального оброблення тексту, завдяки чому будь-яке програмне забезпечення загального призначення, яке має змогу працювати із SGML, може опрацьовувати TEI-сумісні тексти.

У процесі програмного оброблення текстів для формування підкорпусу художньої літератури НКУМ використовуємо адаптований набір тегів (кодів), співвідносний зі стандартним базовим набором тегів системи кодування TEI Lite для прозових творів.

Так, обов'язковими елементами для всіх текстів, поданих у форматі TEI, є *електронний заголовок і власне текст*.

### 1. Електронний заголовок.

Цей елемент тексту TEI містить загальну інформацію про текст (ми обрали й уніфікували систему кодів, яка подає бібліографічні дані та параметри тексту). Заголовок уводиться за допомогою елемента **<headerText>** і має чотири основних компоненти з відповідними атрибутами та значеннями:

<b>&lt;author&gt;</b>	містить повне ім'я автора;
<b>&lt;title&gt;</b>	містить назву твору;
<b>&lt;soursDesk&gt;</b>	групує бібліографічний та жанрово-стильовий опис опрацьовуваного тексту, роблячи це за допомогою таких елементів:
<b>&lt;edition&gt;</b>	подає особливості цієї редакції тексту;
<b>&lt;vol&gt;</b>	містить інформацію про номер періодичного видання, серії, тому тощо; подається за допомогою атрибута <i>type</i> , який має такі значення: <i>number</i> (номер), <i>series</i> (серія), <i>volume</i> (том);
<b>&lt;style&gt;</b>	подає інформацію про стиль художнього твору, що вводиться в корпус;
<b>&lt;genre&gt;</b>	подає інформацію про жанр художнього твору, що вводиться в корпус;
<b>&lt;extend&gt;</b>	подає інформацію про розмір електронного тексту; для зручності розмір вказується в кількості слів (оскільки слово є одиницею виміру самого корпусу), за допомогою атрибута <i>type</i> <sup>1</sup> , який має значення <i>w</i> – скороченої форми терміна <i>word</i> (слово);
<b>&lt;publicationStmt&gt;</b>	групує інформацію про

<sup>1</sup> Глобальний атрибут *type* конкретизує текст через присвоєння йому різних значень і активно використовується у випадках типологічного розмаїття

---

---

	публікації / видання / перевидання опрацьованого тексту;
<publisher>	містить назву видавництва;
<pubPlace>	подає назву місця, де розташоване видавництво;
<date>	містить дату виходу видання (дата може бути подана в будь-якому форматі);
<address>	містить адресу видавництва, представництва, електронної бібліотеки, сайту тощо, тобто місця, звідки взято текст. Передбачені такі значення для уточнення типу адреси: <i>e-mail</i> – для електронної, <i>a-post</i> – для поштової.

## 2. Основна частина – власне текст.

Систему тегів для розмітки основної частини умовно можна поділити на дві групи. Першу групу становлять коди для розділення тексту, виокремлення як основної частини тексту, так і її структурних елементів і субелементів, а саме частин (розділів, книг тощо), абзаців, речень, а також пропусків у тексті. Другу групу формує система тегів для виділення окремих елементів тексту.

Текстову частину вводять за допомогою елемента <text>. Розділення тексту на першому рівні супроводжується використанням таких елементів:

<front>	містить довільну вступну інформацію, яка розміщена перед основним текстом;
<group>	містить монотекст, множину текстів або груп текстів;
<body>	містить усю основну частину одного монотексту, крім того, що стосується вступної чи завершальної частин;
<back>	містить різні додатки і все, що розташоване після основного тексту.



Основна частина прозового тексту може бути згрупована в глави (розділи, підрозділи, книги тощо), або представлена у вигляді простого набору абзаців. У першому випадку використовуємо елемент **<div>**. Взагалі цей тег застосовуємо у випадку виділення будь-якої чітко вираженої частини, наприклад, діалогу без слів автора або віршованого тексту, із типологічним визначенням такої частини за допомогою відповідних елементів. У другому випадку абзац відмічається тегом **<p>**. У межах основної частини всі речення виділяються елементом **<s>**. Пропуски тексту позначаються тегом **<gap>**.

Таким чином, для розділення тексту на другому рівні використовуються такі елементи:

<b>&lt;div&gt;</b>	містить розділ або чітко виражену частину прозового тексту; обов'язковими для цього елемента є атрибути <i>type</i> і <i>n</i> ; останній означає коротку (але зрозумілу) назву або номер розділу; атрибут <i>type</i> подає загальноприйнятну назву для відповідної категорії частини або вказує на тип частини, яку необхідно виділити; в НКУМ для атрибута <i>type</i> використовуються такі значення:
<b>&lt;chapter&gt;</b>	розділ, глава, частина;
<b>&lt;book&gt;</b>	книга;
<b>&lt;poem&gt;</b>	вірш у прозовому творі; розмітка самого вірша відповідно до стандартів TEI з використанням елементів <b>&lt;lg&gt;</b> (виокремлює групу віршованих рядків, що становлять певну структурну одиницю) і <b>&lt;l&gt;</b> (виокремлює віршований рядок);
<b>&lt;speech&gt;</b>	пряма мова;
<b>&lt;p&gt;</b>	відмічає абзаци прозового тексту;
<b>&lt;s&gt;</b>	відмічає речення прозового тексту;
<b>&lt;gap&gt;</b>	відмічає пропуск у тексті.

Наступний перелік тегів та їхніх значень стосується кодування елементів тексту.

<b>&lt;abbr&gt;</b>	скорочення будь-якого типу;
<b>&lt;address&gt;</b>	будь-яка адреса в тексті;
<b>&lt;cit&gt;</b>	цитата з будь-якого іншого документу;
<b>&lt;date&gt;</b>	дата у будь-якому форматі;
<b>&lt;num&gt;</b>	будь-яке число;
<b>&lt;lg&gt;</b>	група віршованих рядків, що становлять певну структурну одиницю;
<b>&lt;l&gt;</b>	віршований рядок;
<b>&lt;name&gt;</b>	будь-яка власна назва; потребує використання атрибута <i>type</i> з такими значеннями:
<antrop>	антропоніми;
<toponim>	топоніми;
<place>	назви країв, країн, територій
<astrolog>	назви астрономічних об'єктів;
<teonim>	назви божеств;
<zoonim>	клички тварин;
<organisation>	назви організацій;
<literat>	назви художніх творів;
<print>	назви друкованих видань
<b>&lt;term&gt;</b>	термін; у разі вживання терміна іноземною мовою передбачено використання атрибута <i>lang</i> із зазначенням перших літер назви мови англійською мовою, наприклад: <i>engl</i> – англійська, <i>kirg</i> – киргизька, <i>chin</i> – китайська, <i>lat</i> – латина; <i>slav</i> – старослов'янська мова тощо;
<b>&lt;sp&gt;</b>	пряма мова;
<b>&lt;speaker&gt;</b>	мовець прямої мови;

<b>&lt;foreign&gt;</b>	запозичення; у разі вживання іноземною мовою також передбачено використання атрибута <i>lang</i> із зазначенням відповідної мови;
<b>&lt;hi&gt;</b>	відмічає слово чи фразу, які графічно відрізняються від загального тексту; потребує атрибут <i>rend</i> , якому, відповідно, надаються наступні значення:
<b>&lt;italic&gt;</b>	курсив;
<b>&lt;bold&gt;</b>	напівжирний шрифт;
<b>&lt;undo&gt;</b>	підкреслений текст;
<b>&lt;ll&gt;</b>	велика літера ( <i>large letter</i> ).

**Приклад:**

```

<header Text>

  <author> В'ячеслав Сахно </author>
  <title> Циркулятор </title>

  <sourceDesc>
    <edition> Сучасність </edition>
    <vol> 3 </vol>
    <style> художній </style>
    <genre> роман </genre>
    <extend type=w> 5000 </extend>
  </sourceDesc>

  <pablStmt>
    <pubplace> Київ </pubplace>
    <date> 2001 </date>
    <address> </address>
  </pablStmt>

</header Text>

<text>

```

<p>  
<s>Я зрідка бачу сни.</><s>Або просто забуваю.</><s>Тому я заздрю людям, які їх бачать і пам'ятають.</><s>Адже вони мають змогу жити у двох різних світах.</><s>Одне життя вдень - нудне й однобарвне, геть інше вночі - непередбачуване й різнобарвне.</><s>Щоправда, вранці настає неабияке розчарування - настирливо дзеленчить будильник, мушиш збиратися на осоружну працю.</><s>Гадаю, якби я міг добирати собі сни до вподоби, на власний смак, то, можливо, й змирився б із білим днем.</>

&lt;/p&gt;

&lt;p&gt;

<s>Мені люба нічна пора.</><s>Тим-то я не люблю спати, а вставати вранці й поготів.</><s>Тобто я е, за прийнятими мірками, совою, або, радше, пугачем чи сичем.</><s>Слово чоловічого роду мені імпонує більше.</><s>Цілком імовірно, що людина насправді нічна істота, бо ж саме ніч вона обрала для мрій, любовців і збереження породи.</>

&lt;/p&gt;

&lt;gap&gt;

&lt;p&gt;

<s>Власне, втрачати мені було нічого.</><s>По мені ніхто б і не заплакав!</><s>Але ж це справжнісіньке самогубство!</><s>Можливо, в мене ще будуть шанси, не такі ризиковані!..</>

&lt;/p&gt;

&lt;p&gt;

<sp>Твоєму життю ніщо не загрожує!</sp> немов гіпнотизуючи мене, провадив <name type=антроп>Богдан</>. <sp>На твоєму місці я б не відмовлявся від такого шансу.</ <s>Він підвівся.</> <sp>Думай, старий! Думай! Остаточну відповідь <gap></sp>

&lt;/p&gt;

&lt;/text&gt;

**Література**

1. Демська-Кульчицька О. М. Основи Національного корпусу української мови. – К., 2005.
2. Передрієнко В. А. Трьох стилів теорія // Енциклопедія. Українська мова. – К.: Видавництво Українська енциклопедія ім. М. П. Бажана, 2000. – С. 640–641.
3. Standard Generalized Markup Language ISO 8879: Information processing – Text and office systems. – Geneva, 1986.
4. Guidelines for Electronic Text Encoding and Interchange / С. М. Sperberg-McQueen, L. Burnard. – 2001. – <http://www.hcu.ox.ac.uk/TEI/P4X/index.html>.

**Тетяна Цимбалюк-Скопненко**, к. ф. н.\*  
Інститут української мови НАН України (Київ)  
УДК 161.2.81'374

**УКЛИНЮВАННЯ ЯК ОДИН ІЗ ВИДІВ ЛЕКСИКО-  
СЕМАНТИЧНИХ ПЕРЕТВОРЕНЬ ФРАЗЕОЛОГІЗМІВ У  
МОВІ ПЕРЕКЛАДУ МИКОЛИ ЛУКАША**

*У статті проаналізовано специфіку вклинювання як лексико-семантичного перетворення у творчій спадщині видатного українського перекладача М. Лукаша, встановлено особливості інкорпорації елементів різних частин мови до структури традиційно вживаної фразеологічної одиниці в перекладних поетичних і прозових текстах.*

У другій половині ХХ ст. інтенсивність інновацій у фразеологічних системах сучасних слов'янських мов помітно зросла. Ці динамічні процеси мають багаторівневий характер, оскільки культурно-політичні події згаданого періоду справили надзвичайний вплив на всі мови європейського ареалу. Лексико-семантичні перетворення у фразеологічній системі літературної мови передусім полягають у тому, що оказіональні елементи доповнюють або змінюють компонентний склад традиційно вживаної одиниці. Використання вклинювання<sup>1</sup> як способу перетворення фразеологічної одиниці (далі – ФО) свідчить „про

---

\* © Т. Цимбалюк-Скопненко, 2005

<sup>1</sup> Термін запровадив О. Кунін [2, с. 13].