



Використання корпусних технологій у дослідницькій та навчальній діяльності

М.В. Надутенко

www.ulif.org.ua

<http://lcorp.ulif.org.ua/dictua/>

**Результати та надбання
Українського мовно-
інформаційного фонду НАН України
на поприщі лінгвістичних
технологій**

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Відкрито лексикографічний ефект в інформаційних системах.
- Побудовано інформаційно-структурну теорію лексикографічних систем.
- Побудовано теорію інтегрованих, індексованих та віртуалізованих Л-систем.
- Розроблено концепцію віртуальних систем професійної взаємодії (спільно з Інститутом кібернетики НАНУ), теорію та технологію віртуальних лексикографічних лабораторій.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено парадигматичну словозмінну Л-систему української мови.
- Виділено близько 3000 парадигматичних класів української лексики.
- Створено Граматичну лексикографічну базу даних обсягом понад 570 тисяч одиниць та відповідну Віртуальну лексикографічну лабораторію (електронний Граматичний словник української мови).
- Розроблено алгоритми побудови фонематичної транскрипції для української лексики та побудовано базу транскрибованої словозмінної парадигми на масиві електронного Граматичного словника української мови.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено парадигматичну словозмінну Л-систему російської мови.
- Створено граматичну лексикографічну базу даних обсягом понад 180 тисяч одиниць та відповідну Віртуальну лексикографічну лабораторію.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Побудовано теорію іменної словозміни турецької мови.
- Кваліфіковано іменну словозмінну парадигму турецької мови обсягом 1808 граматичних значень.
- Досліджено фонетичні процеси та розроблено алгоритми побудови повної словозмінної парадигми турецького іменника.
- Створено ВЛЛ «Грамматичний словник турецької мови» та «Тлумачний словник турецької мови».

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Німецька мова. Розроблено парадигматичну словозмінну Л-систему.
- Іспанська мова. Розроблено парадигматичну словозмінну Л-систему.
- Створено теорію семантичних станів мовних одиниць, яка надає концептуальний апарат для опису лінгвістичних неоднозначностей різних типів та інтеграції граматичного та лексикографічного типів опису мовної системи.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено теорію концептографічних систем. Реалізовано системну концептографію Святого Письма (Євангелія).
- Розроблено систему лінгвістичних алгоритмів автоматичного аналізу українських текстів, яка включає: передморфологічний, морфологічний, синтаксичний й частково семантичний аналіз.
- Побудовано теорію лексикографічної системи українського дієслова та іменника.
- Розроблено лінгвістичну теорію еквівалентів слова. На основі цієї теорії створено цілу низку словників еквівалентів слова.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено теорію лексикографічної системи фундаментального тлумачного Словника української мови.
- Розроблено теорію лексикографічних середовищ, яка стала теоретичним підґрунтям для побудови інтегрованих лексикографічних систем.
- За допомогою цієї теорії було розроблено **Інтегровану лексикографічну систему «Словники України»** - найбільший і найдовершеніший український електронний словник.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено теорію лінгвістичного корпусу української мови. Сформовано УНЛК обсягом близько 180 млн. слововживань.
- Розроблено інтегровану систему двомовної слов'янської лексикографії «МОНДІЛЕКС». Створено 30 двомовних ВЛЛ (мови болгарська, польська, російська, словацька, словенська, українська).
- Розроблено теорію лінгвістичних онтологій.
- Створено Український лінгвістичний портал.

Результати та надбання Українського мовно-інформаційного фонду НАН України

- Розроблено теорію та технологію віртуалізації лексикографічних систем. Здійснено реалізацію УНЛК у хмаринному середовищі та на суперкомп'ютері SKIT-3.
- Сформульовано та обґрунтовано концепцію системи Національних лінгвістичних ресурсів та Національної лінгвістичної інфраструктури.

Електронні ресурси Національної
словникової бази (Віртуальні
лексикографічні лабораторії,
лінгвістичні корпуси, словники,
платформи «LEX» та «ULISS»)



РЕСУРСИ на сайті Українського мовно-інформаційного фонду НАН України

- «Словники України on-line»
 - Тлумачний словник української мови у 20-ти томах (попередня версія)
 - Тлумачний словник української мови у 20-ти томах *
 - Система лінгвістичної взаємодії «ВЛЛ»
 - «Український національний лінгвістичний корпус»
-
- Система лінгвістичної взаємодії «ВЛЛ» (нова версія)
 - Словник української мови за редакцією Б.Д. Грінченка *
 - Словарь русского словоизменения
 - Віртуальна лексикографічна система "MONDILEX"
 - Віртуальна лексикографічна система "Turkic"
 - Український біографічний архів
-
- Сайт "Українсько-російсько-англійський словник зі зварювання" (веб версія)
 - Російсько-українсько-англійський словник зі зварювання
-
- Російсько-українсько-англійський словник з механіки
 - Російсько-українсько-англійський словник з механіки (версія 2)
-
- Російсько-українсько-англійський словник з радіоелектроніки
-
- Віртуальна термінографічна лабораторія з фізики
 - Віртуальна термінографічна лабораторія з біології
 - Віртуальна лексикографічна лабораторія «Українсько-кримськотатарський словник»
-
- "Словники України" (Версія з диску 4.0 (мережева))

- ВЛЛ «Словник української мови»
- ВЛЛ «Етимологічний словник української мови»
- УКРАЇНСЬКИЙ НАЦІОНАЛЬНИЙ ЛІНГВІСТИЧНИЙ КОРПУС
- ВЛЛ «Словарь русского языка»
- ВЛЛ «Словник турецької мови»
- ВЛЛ «Граматичний словник української мови»
- ВЛЛ «Граматичний словник російської мови»
- ВЛЛ «Граматичний словник турецької мови»
- ВЛЛ «Граматичний словник німецької мови»
- ВЛЛ «Граматичний словник іспанської мови»
- ВЛЛ «Словник синонімів української мови»
- ВЛЛ «Українська ономастика»

- 30 ВЛЛ «MONDILEX» для укладання двомовних словників у системі мов: болгарська, польська, російська, словацька, словенська, українська.
- ВЛЛ «Українсько-кримськотатарський словник»
- ВТЛ «ЗВАРЮВАННЯ»
- ВТЛ «ФІЗИКА»
- Інтегрована лексикографічна система «Словники України»
- Українсько-російський та Російсько-український словник зі зварювання
- Українсько-російсько-англійський словник зі зварювання
- Українсько-російсько-англійський термінологічний словник-довідник «Зварювання»
- Російсько-українсько-англійський словник з механіки
- Російсько-українсько-англійський словник з радіоелектроніки
- Шестимовний словник металургійних термінів (українсько-грузинсько-російсько-англійсько-французько-німецький)

Для сучасних досліджень в галузі прикладної лінгвістики, Big Data та інформаційних потоків

необхідним є такий інструмент як

Лінгвістичний Корпус.

Зокрема, Лінгвістичний Корпус Google`а, на якому проводяться основні дослідження, має обсяг приблизно 15 трильйонів слововживань.

ЩО ТАКЕ КОРПУС?

Лінгвістичний корпус – інформаційно-довідкова система, яка складається із певних текстів в електронній формі. Це безкоштовний сайт, створений групою вчених Українського мовно-інформаційного фонду НАН України та Національного центру Мала академія наук. Корпус постійно поповнюється й розвивається.



Лінгвістичні корпуси

- Корпус Google:

Понад 1 000 000 000 000 000 000 слівовживань

- **Національні корпуси – декілька сотень мільйонів (до 10^9) слівовживань**
- **Український національний лінгвістичний корпус – 180 млн. слівовживань (віртуалізована платформа ULISS)**

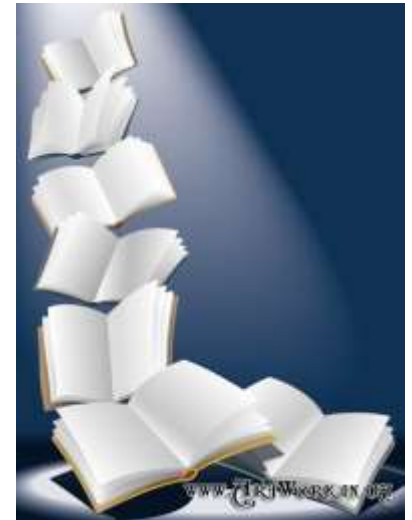
Особливості корпусу:

Корпус має дві важливі особливості, які відрізняють його від інших електронних бібліотек :

1. Репрезентативність та збалансованість складу текстів -

корпус містить практично всі типи текстів:

- класична та сучасна художня література - проза й поезія;
- словники (близько 240);
- мемуарно-біографічна література;
- журнальна публіцистика й літературна критика;
- газетна публіцистика й новини;
- наукові, науково-популярні й навчальні тексти;
- релігійні й релігійно-філософські тексти;
- виробничо-технічні, офіційно-ділові і юридичні тексти.



2. Має розмітку. Розмітка надає особливу додаткову інформацію про властивості текстів, які входять до корпусу. Професійно зроблена розмітка дозволяє швидко та ефективно знайти ті слова, форми та конструкції, які потрібні досліднику.

Як можна використовувати лінгвістичний корпус?



- підготовка навчальних матеріалів для занять: розроблення вправ, тестів та контрольних;

- організація самостійної роботи учнів;
- організація науково-дослідної роботи.



Як можна використовувати лінгвістичний корпус?

- - наводити приклади з корпусу в роздаткових матеріалах;
- - продемонструвати пошукові можливості корпусу;
- - роздавати приклади для інтерпретації;
- - домашнє завдання - провести пошук у корпусі й представити зібраний матеріал;
- - наукові дослідження учнів – визначити тему, зібрати в корпусі матеріал й інтерпретувати його.

- Можна вибрати короткі приклади;
- Користувачі розуміють, що вони можуть обрати тему, зібрати матеріал і інтерпретувати його, - так вони стануть самостійними вченими
- Можна проілюструвати тими або іншими прикладами вживання слів у контексті

Які можливості надає корпус?

- пошук текстів;
- пошук додаткового матеріалу;
- довідково-інформаційна робота;
- дослідження творчості окремих авторів;
- перевірка правильності вживання слів;
- матеріали для науково-дослідної роботи.



Лінгвістичний корпус



Український мовно-інформаційний фонд НАН України
УКРАЇНСЬКИЙ НАЦІОНАЛЬНИЙ ЛІНГВІСТИЧНИЙ КОРПУС

Name: УКРАЇНСЬКИЙ НАЦІОНАЛЬНИЙ ЛІНГВІСТИЧНИЙ КОРПУС

Version: 5.4.24.1

Publisher: Український мовно-інформаційний фонд НАН України

Run

Лінгвістичний корпус

Реєстрація користувача



ULISS "Український національний лінгвістичний корпус"

Корпус

Український національний лінгвістичний корпус

Мова інтерфейсу

Ukrainian

Користувач

Пароль

Так

Скасувати

© Український мовно-інформаційний фонд НАН України, 2005-2014



Лінгвістичний корпус

ULISS "Український національний лінгвістичний корпус"

Адміністрування Бібліотека ОНТОГРАФІВ (Система ТОДОС)

пошук за б/о повнотекстовий пошук

Відкрити текст Додати поточний документ

Характеристика

- Назва видання
- Прізвище та ініціали
- Стиль
- Серія
- Жанр
- Ключові слова
- Видавництво
- Місце видання
- Мова видання
- Рік видання (3)
- Рік видання (По)

Пошукова умова

Рекомендована література Так

Результат пошуку - 1859

№	Результат пошуку
10512	Абхазькі народні казки : Казка / перекл. Іщенко Є.. – http://chytanka.com.ua/static/753.ukr.html
7134	Авраменко О. Заборонені чари : Повість. – http://vesna.org.ua/txt/avramenko/zaboroneni_czary.html
10040	Авраменко О. Небо, повне зірок : Повість. – http://javalibre.com.ua/java-book/book/2925946
10038	Авраменко О. Первісна. Дорога на Тір Мінеган : Роман. – http://javalibre.com.ua/java-book/book/2925946
10039	Авраменко О. Первісна. У вирі пророцтв : Роман. – http://javalibre.com.ua/java-book/book/2925946

Імпорт

Бі...

Абхазькі народні
<http://chytanka.com>

Авраменко О. За
vesna.org.ua/txt/

Авраменко О. Не
javalibre.com.ua/

Авраменко О. П
Роман. – <http://ja>
book/2925946

Авраменко О. П
<http://javalibre.com>

Авраменко О. Я,
www.ukrcenter.co

Авраменко О., А
долонях. Жменя
vesna.org.ua/txt/
avramenko_zori_y

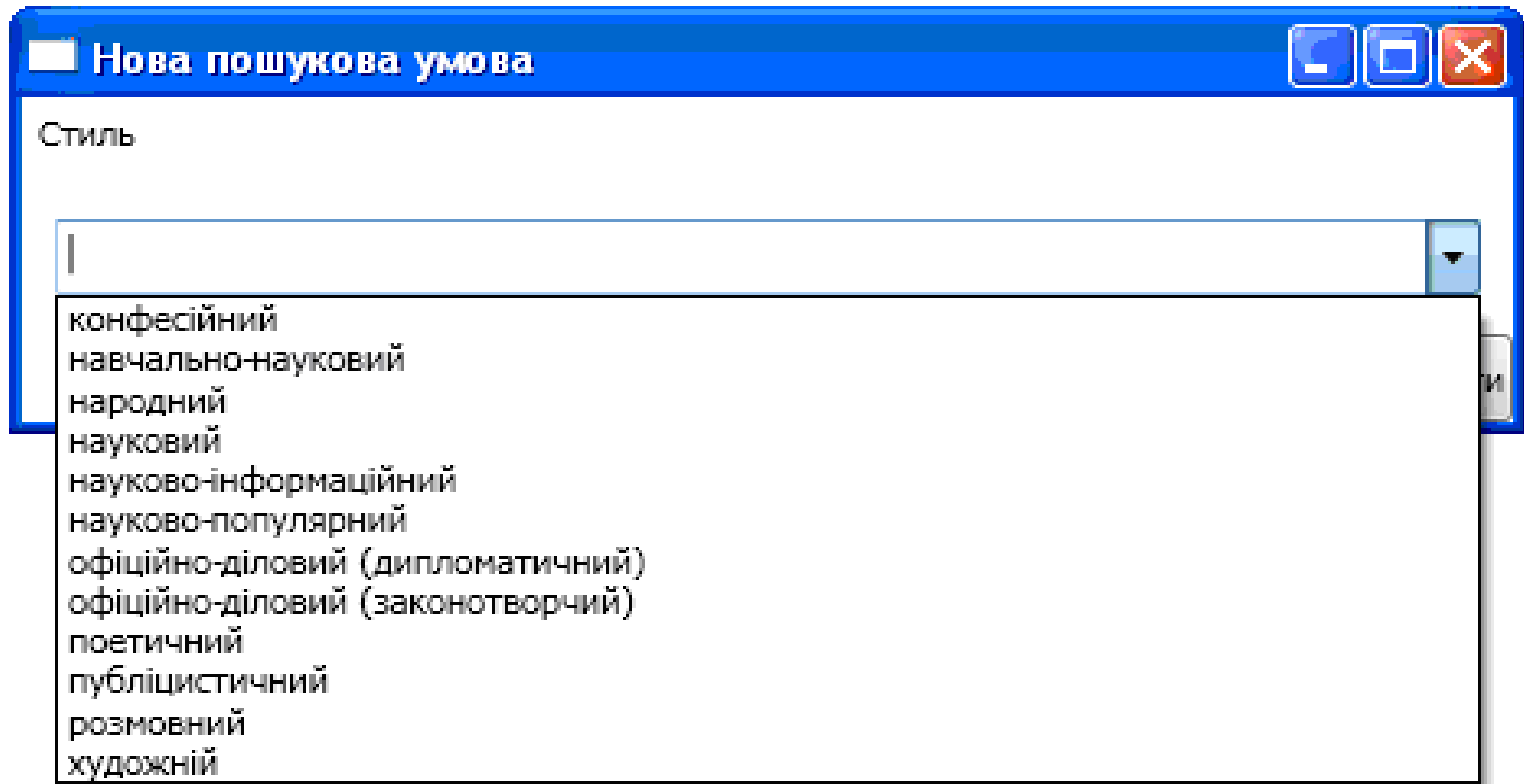
Лінгвістичний корпус: Як здійснюється пошук?

За бібліографічною ознакою:

- Назва видання
- Прізвище та ініціали
- Стил
- Серія
- Жанр
- Ключові слова
- Видавництво
- Місце видання
- Мова видання
- Рік видання (З)
- Рік видання (По)

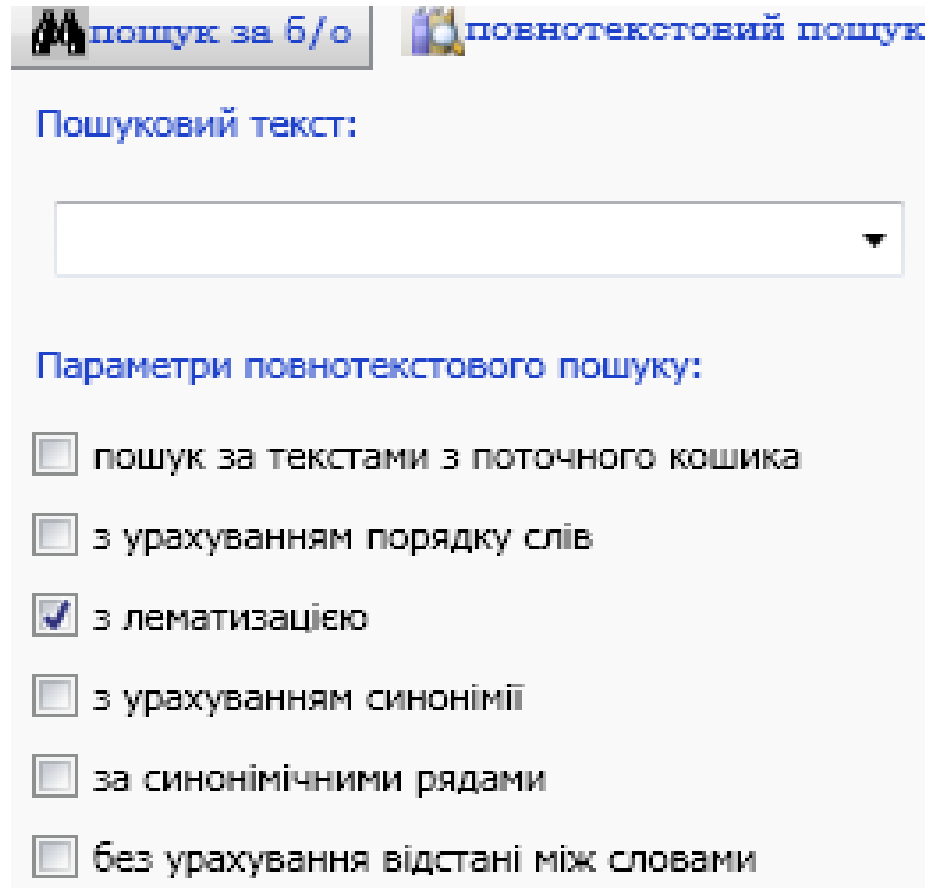
Лінгвістичний корпус: Як здійснюється пошук?


Для кожної ознаки передбачені свої пошукові умови :




Лінгвістичний корпус: Як здійснюється пошук?

Повнотекстовий пошук слова або словосполучення за параметрами:



 пошук за б/о

 повнотекстовий пошук

Пошуковий текст:

Параметри повнотекстового пошуку:

- пошук за текстами з поточного кошика
- з урахуванням порядку слів
- з лематизацією
- з урахуванням синонімії
- за синонімічними рядами
- без урахування відстані між словами

Приклад завдання з використанням корпусу

Синоніми, пароніми у діловому мовленні

Пароніми — слова близькі за звуковим складом і звучанням, та різні за значенням. Різниця може бути в одній букві.

Слова «компанія» — «кампанія».

Приклади використання цих слів у текстах

Результати пошуку слова «компанія» з урахуванням синонімії

1171	Голубець М. Велика історія України : Монографія - Л., 1993.	
245	Гомас О.М. Варіантність і синонімія в словосполученні й реченні : Автореф. дис.. канд. філол. наук - К., 2000.	
5792	Гончар О.Т. Далекі вогнища : Повість	
578	Гончар О.Т. Твори в дванадцяти томах - К. : Наук. думка, 2001. - ISBN 966-00-0710-8	
1007	Горбатюк Н.П. Роль української діаспори в формуванні українсько-канадських відносин : Автореф. дис.. канд. іст. наук - О., 1999.	
828	Горлієнко В.В. Православні конфесії в Україні періоду другої світової війни (вересень 1939 – вересень 1945 рр.) : Автореф. дис.. канд. іст. наук - К., 1999.	

Знайдено: 3050 Час пошуку: 0,203125 Розмір контексту: 500

[Голубець М. Велика історія України : Монографія - Л., 1993.](#)

Контексти

Контекст № 446
Контекст № 447
Контекст № 448
Контекст № 449
Контекст № 450
Контекст № 451
Контекст № 452
Контекст № 453
Контекст № 454
Контекст № 455
Контекст № 456
Контекст № 457
Контекст № 458
Контекст № 459
Контекст № 460
Контекст № 461
Контекст № 462
Контекст № 463
Контекст № 464
Контекст № 465
Контекст № 466
Контекст № 467

нова хвиля кочовиків знову
загрозила Європі, скошуючи свій перший розмах і лють на культурі української
землі. Пропашою називає пок. історик С. Томашівський цю добу в історії України й
східньої Європи, бо вона знищила «дотеперішні початки громадської організації,
віками вироблену культуру й завершила край у стан первісного варварства».
А всеж таки... На велетенській території, над якою стільки століть шаліли бурі
від сходу, на якій пролито стільки людської крові й знищено стільки людської
праці,

Результати пошуку слова «кампанія» з урахуванням синонімії

1790	Білоус Д. Чари Барвінкові : Поезії / Бібліотека з української літератури для школярів - К. : Студія "Негоціант", 2000. - (Золота скарбниця України). - ISBN 966-7423-50-6	
3337	Бондарчук Ю.П. Залізничний транспорт України в умовах утвердження адміністративно-командних методів управління народним господарством (кін. 20-х — 30-ті рр.) : Автореф. дис.. канд. іст. наук / Дніпропетровський державний університет - К., 1999.	
3334	Борисенко В.М. Літературні організації в суспільно-політичному житті України (1920-1932рр.) : Автореф. дис.. канд. іст. наук / Київський університет імені Тараса Шевченка - К., 1999.	
1539	Боротьба з організованою злочинністю і корупцією (теорія і практика)т.1 : Журнал; Науково-практичний журнал / Міжвідомчий науково-дослідний центр - К., 2000.	
1538	Боротьба з організованою злочинністю і корупцією (теорія і практика)т.2 : Науково-практичний журнал / гол. ред. Кондратьєв Я.Ю.; Міжвідомчий науково-дослідний центр - К., 2000.	
1541	Боротьба з організованою злочинністю і корупцією (теорія і практика)т.4 : Науково-практичний журнал / Міжвідомчий науково-дослідний центр - К., 2001.	

Знайдено: 704 Час пошуку: 0,015625 Розмір контексту: 500

Білоус Д. Чари Барвінкові : Поезії / Бібліотека з української літератури для школярів - К. : Студія "Негоціант", 2000. - (Золота скарбниця України). - ISBN 966-7423-50-6

Контексти

Контекст № 1
Контекст № 2
Контекст № 3
Контекст № 4

Мовляв, щоб розпізнати,
що втнув якийсь чудака,
слід, кажуть, підписати:
«Се лев, а не собака»!

КАМΠΑНІЯ І КОМПАНІЯ

Запитала Надя на уроці мови:
— Юрію Петровичу, чом це так бува:
як пишу кампанія — правите на о ви,

Перегляд повного тексту

■ Перегляд повного тексту

«Сиди, квітку твою, підписати».

«Се лев, а не собака»!

КАМΠΑНІЯ І КОМПАНІЯ

Запитала Надя на уроці мови:

— Юрію Петровичу, чом це так бува:
як пишу кампанія — правите на о ви,
а пишу компанія - правите на а?

Осміхнувся учитель і сказав Надії:

— Придивись, подумай, — різні це слова:
мовиться про заходи чи воєнні дії —
пишеться кампанія — з літерою а.

А коли це група, де зібрались люди
на пікнік, на гулі чи на торжество —
написання слова зовсім інше буде —
пишеться компанія — з літерою о.

Аж Данило раптом піднімає руку,
мав завжди за муку цю складну науку:

— А якби полегшить, — пропонує Даня, —
та ввести єдине слів цих написання?

Середовище ТОДОС

Середовище ТОДОС (Трансдисциплінарні Онтологічні Діалоги Об'єктно-орієнтованих Систем) - створене з використанням корпусних технологій та лінгвістичних онтологій

Трансдисциплінарність – спосіб розширення наукового світогляду, що полягає в розгляді того чи іншого явища поза рамками будь-якої однієї наукової дисципліни. За допомогою Т. здійснюється класифікація і систематизація формального взаємозв'язку окремих дисциплінарних знань.

Онтологія - є відображенням певної теорії й може бути представлена як активна система знань, яка включає в себе множину об'єктів (позначених певними термінами), що пов'язані з описами, а також формальні аксіоми, які обмежують інтерпретацію і спільне вживання цих термінів.

Діалог - форма і організація спілкування, комунікації, обміну висловлюваннями.

Об'єктно-орієнтований підхід - спеціальним способом організовані множини, об'єкти яких мають загальні властивості.

Середовище ТОДОС

- забезпечує можливість для автоматизованого подання знань у вигляді динамічної онтології – сукупності термінів дисципліни (чи будь-якої їх кількості) та зв'язків між ними.

ВІДЕО

ДЯКУЮ ЗА УВАГУ!