

## КОЛЬОРОВА МАПА ПОЕТИЧНОГО СПАДКУ ВАСИЛЯ СТУСА У MATHEMATICA

**Реферат.** У статті описано ідею та реалізацію процесу створення кольорової мапи (КМ) тексту на основі поетичного спадку Василя Стуса. КМ – це множина кольорових квадратиків (або інших фігур), кожен із яких представляє конкретну кольороназву в оригінальному тексті. Цілком об'єктивний результат унаочнює розподіл у тексті конкретних прикметників на позначення кольору. Оригінал статті написано у форматі *CompuTable Document Format (CDF)* – обчислюваний документ, і може бути використаний для довільного тексту українською мовою, з урахуванням особливостей відмінювання і навіть словотвору.

**Ключові слова:** кольороназва, мовна модель, відмінювання.

Сьогодні є чимало інструментів для опрацювання природного мовлення, зокрема, з високим ступенем автоматизації. Сучасні мови програмування дозволяють ґрунтовно обробляти текстові дані на рівні окремих символів, групи символів (фраз і речень), у перспективі мають з'явитися функції опрацювання семантики. Кількість вбудованих рядкових функцій зростає, і це розширює потенціал розробника і допомагає заощаджувати час, необхідний для власноручного створення таких функцій і процедур. Наприклад, процедуру визначення різниці між двома рядками – так званої мінімальної відстані редагування або числа Левенштейна – розробник на Python має написати вручну [Jurafsky 2009]. Втім, у деяких сучасних систем ця процедура входить до вже вбудованих. Обчислювальне програмне середовище *Mathematica* від Wolfram Research, яке широко використовується в багатьох наукових, інженерних, математичних і обчислювальних проектах [Wellin 2013], а нами було застосоване у галузі лінгвістики, наприклад, має вбудовану функцію *EditDistance [X, Y]*. На нашу думку, основні алгоритми автоматичного опрацювання природного мовлення, як-от лематизація, синтез словоформ, визначення граматичних класів і підкласів, синтаксичного аналізу або навіть автоматизованого перекладу, які описано в [Баранов 2003; Волошин 2004; Дарчук 2008; Карпіловська 2006; Марчук, 2000; Партико 2008], згодом стануть вбудованими процедури.

Ця стаття в оригіналі написана у форматі CDF – обчислюваний документ для системи *Mathematica*, або *Notebook (\*.nb)* – і її можна завантажувати [Данилюк 2014], оскільки код у друкованій версії буде наведено тільки частково. Для перегляду CDF вам знадобиться безкоштовна програма від [<http://www.wolfram.com/cdf-player/>].

Отже, головна мета статті полягає в описі процесу, інструментів і безпосередньо коду для автоматичного генерування кольорової мапи для довільного тексту українською мовою, і зокрема, для поетичного спадку Василя Стуса. Кольорова мапа (КМ) – це множина кольорових квадратиків (або інших фігур), кожен із яких представляє конкретну кольороназву в оригінальному тексті. Кожне використання прикметника на позначення кольору – «білий», «чорний», "червоний", «золотий» – у КМ буде представлено квадратиком відповідного кольору. Отже, можна отримати повне й абсолютно об'єктивне представлення лексики конкретного тексту і певні риси "картини світу", концепти окремих кольорів у літературних творах. Це питання є досить актуальним в українській лінгвістиці і філології, тому інструмент для автоматизованого пошуку кольороназв у довільному тексті, на нашу думку, є необхідним.

КМ можна вважати автоматизованою інфографікою для візуалізації мовленнєвих даних. Загальна ідея належить Тетяні Дружняєвій із видання «Esquire», окремі елементи коду

запропоновано Романом Осиповим (Московський державний університет тонких хімічних технологій).

Ми поділили дослідження на **декілька завдань**: 1) отримати всю можливу статистичну інформацію з тексту для подальшого аналізу; 2) створити процедури (мовою для *Mathematica*), щоб працювати з окремими словами та реченнями; 3) побудувати мовну модель для називання кольору з урахуванням відмінювання й окремих випадків словотвору; 4) використати модель для генерування КМ творчого спадку Василя Стуса, і описати перспективи.

Об'єкт дослідження – текст збірок Василя Стуса («Круговерть» (1965), «Зимові дерева» (1970), «Веселий цвинтар» (1971), «Час творчості / Dichtenszeit» (1972), «Палімпсести»). Формально це текстовий файл формату txt, підготовлений для обробки (в Unicode, кожен токен відділений пробілом). Предметом дослідження є використання кольороназв, представлених у вигляді КМ.

Розпочнемо з першим завданням. Текстовий файл – stus.txt – має бути в тій самій папці, що й робочий файл \*.nb, і ми читаємо дані з нього у змінну stus:

```
Short[stus=Import[FileNameJoin[{NotebookDirectory[],"stus.txt"}]]]
```

Коли ми звертаємося до stus, то працюємо з усім текстом і можемо отримати основні статистичні дані – скільки символів він містить:

```
StringLength[stus]
```

```
579771
```

або скільки рядків:

```
StringCount[stus,"\n"]
```

```
23293
```

або скільки приблизно речень (за умови, що речення може закінчуватися крапкою, знаком оклику або знаком питання):

```
StringCount[stus,{».»»?»!»!»}]
```

```
10689
```

або які символи використовуються в тексті, у відсортованому вигляді:

```
Union[Characters[stus]]
```

```
{!, *, (, ), _, -, ` , [ , ] , < , > , . , , , ; , " , ? , ' , / , : ,  
, , [e, [I, [i, А, Б, В, Г, Д, Е, Ж, З, И, Й, К, Л, М, Н, О, П, Р, С, Т, У, Ф, Х, Ц, Ч, Ш,  
Щ, Э, Ю, Я, а, б, в, г, д, е, ж, з, и, й, к, л, м, н, о, п, р, с, т, у, ф, х, ц, ч, ш, щ, ъ,  
ы, ь, э, ю, я, [e, [i, [i, [Г, [Г, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, а, А, б, В, С, d, D, e,  
g, i, I, j, k, l, m, M, n, N, o, p, P, r, s, T, u, U, V, w, X, y, z, -, ..., «, -, », ", ', "}
```

Зверніть увагу на те місце, де опинилися після сортування деякі специфічні українські літери (їх виділено). Причина в тому, що в стандартній таблиці символів вони розташовані не у загальному списку кирилиці, а на випадкових позиціях. Цю незручність необхідно враховувати у разі використання регулярних виразів (РВ) на кшталт "всі українські літери від А до Я". На практиці РВ  $[/C-Г]+/$  досить добре для підходить для запиту "всі українські літери", але він повертає кілька не-українських символів – Э, ъ, ы, э.

Наступним кроком переходимо до обробки слів. Змінна *allWords* містить всі словоформи зі змінної *stus* без розрізнення великих і малих літер – через заміну за допомогою РВ:

```
Short[allWords = Sort[Tally[DeleteCases[StringSplit[StringReplace[StringReplace[stus,  
Thread[Join[CharacterRange["A", "Я"], {""Є", "І", "Ї", "Г"}] -> Join[CharacterRange["a", "я"], {""є",  
"і", "ї", "г"}]]], RegularExpression["[^" <> StringJoin@Join[CharacterRange["a", "я"], {""є", "і",  
"ї", "г"}] <> "]" ] -> " ", " ", ""], #1[[2]] > #2[[2]] &], 5]
```

Кількість токенів у *stus* сягає 93442, а унікальних словоформ – 23182. Ось 100 найчастотніших:

```
allWords[[1";;" 200]][[";";",1]]
{"і", "не", "в", "на", "як", "ти", "що", "а", "у", "з", "я", "за", "до", "й", "де", "мов", "бо",
"ні", "мене", "так", "вже", "тебе", "ще", "ж", "тільки", "то", "все", "по", "чи", "коли", "із", "це",
"та", "світ", "мій", "себе", "од", "мені", "о", "життя", "тобі", "нас", "над", "щоб", "він", "під",
"твій", "ніби", "хай", "ми", "там", "від", "але", "моя", "ані", "серце", "душу", "є", "немов",
"без", "б", "тут", "смерть", "аж", "лиш", "між", "день", "душі", "той", "аби", "неначе", "всі",
"очі", "тобою", "край", "нам", "небо", "ніч", "геть", "сон", "те", "нема", "хоч", "уже", "пам",
"про", "буде", "десь", "душа", "ось", "вона", "наче", "свій", "його", "хто", "серця", "для",
"твоя", "цей", "руки"}
```

Для створення можливості шукати конкретну словоформу в тексті будемо функцію *wordPosition*, яка повертає масив символів позицій.

```
replacements=Thread[Join[CharacterRange["a","я"],{"e","i","ї","r"}]-
>Join[CharacterRange["A","Я"],{"C","Г","І","Г"}]]
```

Ось результат для слова *тополя*:

```
wordPosition["тополя"]
{{219414, 219421}, {408610, 408617}, {426728, 426735}, {484326, 484333}}
```

Якщо ми знайдемо позиції для крапок (та інших символів у кінці речення – “!”, “?”, “...”), можна буде отримати ціле речення (послідовність символів між двома розділовими знаками) і записати у змінної *sentence*.

```
Short[dots = #[[1]] & /@ StringPosition[stus, {".", "?", "!", "[Ellipsis]"}, 5];
sentence[{min_, max_}]:=Block[{start=Select[Nearest[dots, min, 10], #<min&][[1]]+1, end=Select[Nearest[dots, max, 10], #>max&][[1]]}, StringTake[stus, {start, end}]]
```

Поєднання *wordPosition* і *sentence* виведе конкорданс для конкретної словоформи:

```
Grid[Transpose@{StringReplace[sentence /@ wordPosition["тополя"],
"\n" -> ""}], Background -> {None, {{Orange, LightGray}}},
ItemStyle -> Directive[16, Bold, FontFamily -> "Arial"],
Alignment -> Left, Dividers -> All]
```

<b>І вже дзвінка тополя виростає душі окраденої.</b>
<b>Тополя ламле руки — їй сил нема — пірвати тіло в лет.</b>
<b>До смертного дрозубачу — тополя до мене спішить.</b>
<b>Перелетіть мене, перелетіть через дроти, паркани і горожі, о, Україно, до смертного дрозубачу тополя твоя шелестить.</b>

Тепер треба визначити загальну колірну модель для побудови КМ. Її релевантність і глибина для української мови – зокрема, тісність лематизації і докладність – мають вирішальне значення для якості КМ. Перший крок – знайти прикметники, які прямо позначають кольори і складаються з однієї словоформи. Заносимо ці прикметники у масив *tc* (таблиця кольорів) і описуємо їх за моделлю RGB (фрагмент коду досить довгий і підходить тільки для перегляду в цифровій версії). Це прикметники: *білий, червоний, зелений, синій, жовтий, чорний, сірий, рожевий, коричневий, блакитний, пурпурний, пурпуровий, оранжевий, помаранчевий, фіолетовий, амарантовий, буришинний, аметистовий, абрикосовий, аквамаринний, арсеновий, спаржевий, бежевий, латунний, бронзовий, брунатний, кармінний, морквяний, лазуровий, каштановий, шоколадний, цинамоновий, кобальтовий, мідний, кораловий, кукурудзяний, блаватний, кремовий, малиновий, джинсовий, смарагдовий, баклажановий, ляльний, золотий, індиго, нефритовий, хакі, лавандний, лимонний, бузковий, малахітовий, гірчичний, оливковий, помаранчевий, ліловий, персиковий, грушевий, барвінковий, сливовий, бурий, іржавий, шафрановий, сапфіровий, багряний, срібний, болотний, мандариновий, будяковий, бірюзовий, ультрамаринний, фіолетовий, пшеничний.*

А ось фрагмент представлення *tc*:

білий	GrayLevel[1]
червоний	RGBColor[1, 0, 0]
зелений	RGBColor[0, 1, 0]
синій	RGBColor[0, 0, 1]
жовтий	RGBColor[1, 1, 0]
чорний	
сірий	GrayLevel[0.5]
рожевий	RGBColor[1, 0.5, 0.5]
коричневий	RGBColor[0.6, 0.4, 0.2]
блакитний	RGBColor[0, 1, 1]
пурпурний, пурпуровий	RGBColor[1, 0, 1]
оранжевий, помаранчевий	RGBColor[1, 0.5, 0]
фіолетовий	RGBColor[0.5, 0, 0.5]
амарантовий	RGBColor[ $\frac{229}{255}$ , $\frac{43}{255}$ , $\frac{16}{51}$ ]

Колірна модель також включає в себе правила для побудови різних словоформ і пошуку їхніх лем. Так, *червоний*, *червоного*, *червона*, *червоної* тощо будуть представлені лемою *червоний* і червоним квадратиком. Для цього використовуємо синтез словоформ за словником основ і закінчень (отримання всіх можливих словоформ для кожного прикметника в таблиці кольорів), потім знаходимо позиції для цих словоформ у змінній *stus*, присвоюємо ці позиції конкретній лемі. Нарешті, діапазон лем, як вони з'являються в тексті, замінюємо на кольорові квадратики.

Відмінювання прикметника в українській мові включає тверду (*червоний*) і м'яку групи (*синій*), стягнені (*червона*) і нестягнені форми (*червоная*). Змінні *ColorEdningsTv* і *ColorEdningsMk* містять відповідно закінчення для стягнених і нестягнених форм твердої групи та стягнених і нестягнених форм м'якої групи:

```
ColorEdningsTv={"ий","ого","ому","им","ім","е","а","ої","ій","у","ою","і","их","им",
"ими","еє","ая","ую","ії"};
```

```
ColorEdningsMk={"ій","ього","ьому","ім","є","я","ьої","ю","ьою","і","іх","імі","еє",
"еє","яя","юю","ії"};
```

Змінна *colorRules* містить згенерований масив усіх можливих словоформ для називання кольорів (поєднання основ із *tc* і закінчень із *ColorEdningsTv* і *ColorEdningsMk*) зі вбудованим механізмом врахування афіксів (*білий* – *біленький*). Фрагмент коду:

```
colorRules=Flatten[ {Thread[Flatten[Table[{"біл",
"біленьк"}][[v]]<>ColorEdningsTv[[u]],{v,2},{u,Length[ColorEdningsTv]}]->White],
Thread[Table["син"<>ColorEdningsMk[[u]],{u,Length[ColorEdningsMk]}]->Blue],...
Thread[Table["пшеничн"<>ColorEdningsTv[[u]],{u,Length[ColorEdningsTv]}]-
>RGBColor[245/255,222/255,179/255]]}]
```

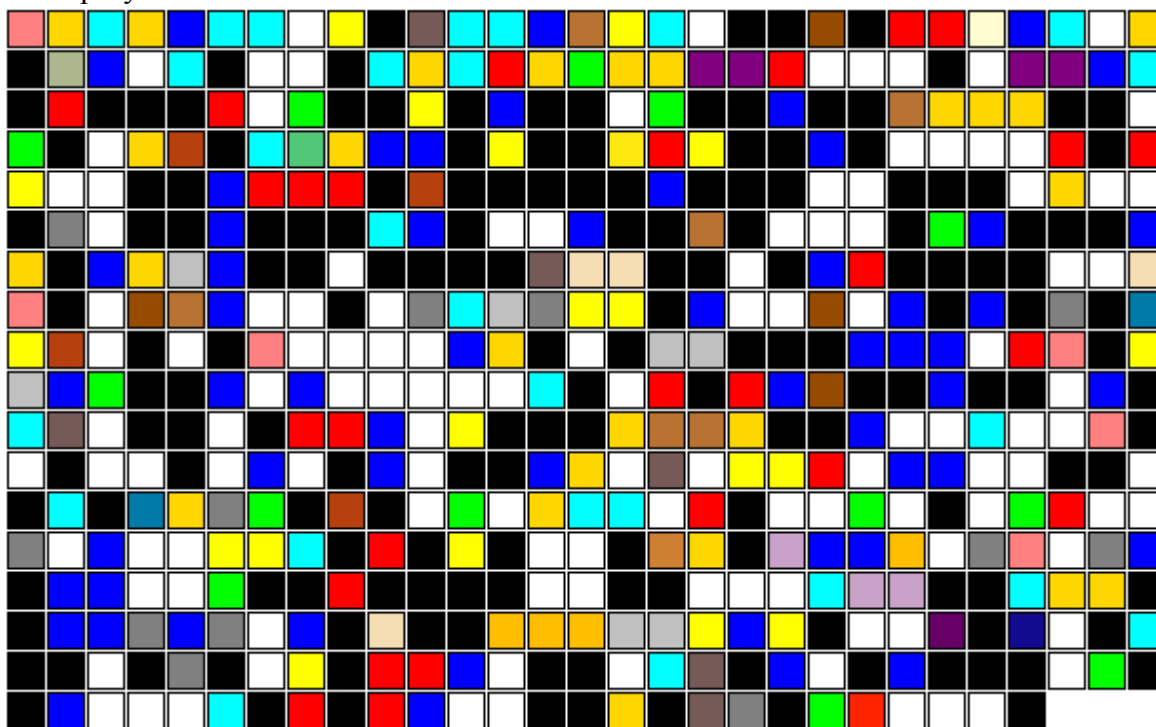
Позиції згенерованих словоформ записуємо у змінній *colorInformationPre*, а потім сортуємо у *colorInformation*:

```
colorInformationPre={#,wordPosition[#]}&/@colorRules[;];
colorInformation=Sort[Flatten[Partition[Riffle#[[2]],#[[1]],{2,-1,2}],2]&/
@colorInformationPre,1],Mean[#1[[1]]]<Mean[#2[[1]]]&]
```

І, нарешті, будуємо КМ:

```
Panel[Grid[Partition[Graphics[{{#,EdgeForm[Black],Rectangle[]},ImageSize-
>20]}&/@(colorInformation[;],2)]/.colorRules),30,30,1,""],Spacings->{0,0}]
```

Ось результат:



У цифровій версії цієї статті можна побачити реальне речення, клікнувши на конкретному квадратику.

Загальний висновок: процедура створення КМ для українського тексту в *Mathematica* є досить складною через виключення окремих символів (*є, и, ї, і*) в "нормального" списку в Unicode, необхідність враховувати, що кириличні літери не включені в список «правильних» символів слів, як-от латиниця або цифри. Ефективно розв'язати цю проблему допомагає створення спеціальних функцій. Сама система *Mathematica* є потужним середовище для досліджень із комп'ютерної лінгвістики з багатою мовою програмування високого рівня, тому її можна рекомендувати для вивчення студентами філологічних відділень в українських вишах.

Текст поетичного спадку Василя Стуса є надзвичайно багатим на кольороназви, основними з яких є чорний, білий, сірий, голубий, синій, червоний, зелений. Порівняння КМ збірок Василя Стуса з КМ інших творів літератури і фольклору може бути темою для подальшого вивчення. Як пише автор дослідження [Ковтун 2009: 51], «колірна ознака з'явилася в мові в діахронічній послідовності. У народній творчості спочатку переважають означення білого та чорного кольорів, за ними йде червоний (тріада білий — чорний — червоний), після нього — зелений і жовтий, далі — синій і брунатний». Виглядає, що творчість Василя Стуса, настільки своєрідна у плані лексики, в аспекті використання кольороназв є близькою до фольклору.

Описаний спосіб побудови КМ на основі позиції лексеми на позначення кольору є базовим. Його можна поліпшити через додавання до колірної моделі похідних прикметників (*білявий, чорнючий*), дієслів (*біліти*) та іменників (*чорнота*), деяких назв вторинних відтінків (*ясно-зелений*), тісно пов'язаної лексики – *вороний, гнідий* тощо.

#### Список використаної літератури

1. Баранов 2003: Баранов, А.Н. Введение в прикладную лингвистику [Текст] / А. Н. Баранов. – М. : Едиториал УРСС, 2003. – 360 с. – ISBN: 5-8360-0196-0.
2. Волошин 2004: Волошин, В. Г. Комп'ютерна лінгвістика : Навчальний посібник [Текст] / В. Г. Волошин. – Суми : ВТД «Університетська книга», 2004. – 382 с. – ISBN: 966-680-134-5.

3. Данилюк 2014: Данилюк, І.Г. Аналіз тексту "Кобзаря" Тараса Шевченка в середовищі Mathematica: символи, слова і кольори [Текст]. – Access mode : URL : <https://app.box.com/stus>. – Title from the screen.
4. Дарчук 2008: Дарчук, Н.П. Комп'ютерна лінгвістика: Автоматичне опрацювання тексту [Текст] / Н.П. КДарчук. – К. : Видавничо-поліграфічний центр «Київський університет», 2008. – 351 с. – ISBN 978-966-439-079-5.
5. Карпіловська 2006: Карпіловська, Є. А. Вступ до прикладної лінгвістики : комп'ютерна лінгвістика. Підручник [Текст] / Є. А. Карпіловська. – Донецьк : ТОВ «Юго-Восток, Лтд», 2006. – 188 с. – ISBN 966-374-078-7.
6. Ковтун 2009: Ковтун, Л. Український колористичний код світотворення // Вісник Київського національного університету імені Тараса Шевченка. Українознавство. – Вип. 13 / КНУ імені Тараса Шевченка. – Київ : ВПЦ "Київський університет", 2009. – ISSN 1728-2330
7. Марчук 2000: Марчук, Ю. Н. Основы компьютерной лингвистики [Текст] / Ю. Н. Марчук. – М. : Народный учитель, 2000. – 320 с. – ISBN 5-17-039480-2.
8. Партико 2008: Партико, З.В. Прикладна і комп'ютерна лінгвістика: Вступ до спеціальності [Текст] / З.В. Партико. – Львів: Афіша, 2008. – 224 с. – ISBN 978-966-325-092-2.
9. Jurafsky 2009: Jurafsky, D., Martin, J. H. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition [Текст]. – Prentice Hall, 2009. – 988 p. – ISBN 0-13-095069-6.
10. Wellin 2013: Wellin, Paul R. Programming with Mathematica, An Introduction. – Cambridge, 2013. – 711 pp. – ISBN: 9781107009462

#### Аннотация

#### **Данилюк І.Г. Цветовая карта поэтического наследия Василия Стуса в MATHEMATICA.**

В статье описаны идея и реализация процесса создания цветовой карты (ЦК) для произвольного текста вообще, и в частности для поэтического наследия Василия Стуса. ЦК - это множество, сетка цветных прямоугольников (или других фигур), каждый из которых относится к определенному названию цвета (лексеме) в исходном тексте. Полностью объективный результат демонстрирует дистрибуцию соответствующих прилагательных. Оригинал статьи создан в новейшем CDF-формате вычисляемого документа и может быть применен к произвольному тексту на украинском языке с учетом словоизменения и словообразования.

*Ключевые слова:* название цвета, языковая модель, словоизменение.

#### Summary

#### **Danyluk I. G. Color map of Vasyl Stus poetry with MATHEMATICA.**

The article describes an idea and its realization process for creating Color Map (CM) for the text – in that case Vasyl Stus poetry. CM is a composition, a grid made of colored rectangles (or any other figures) – one for every name of color in the original text. Absolutely objective result shows visually the distribution of particular adjectives. The original article is in Computable Document Format (CDF) and is suitable for random text in Ukrainian considering inflection and even derivation.

*Keywords:* color name, language model, inflection