

# **Відкриті дані в Україні**

Навчальний посібник

Олександр Краковецький

# Відкриті дані в Україні

## Навчальний посібник

Олександр Краковецький

This book is for sale at <http://leanpub.com/opendatainua>

This version was published on 2017-03-20



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2016 - 2017 Олександр Краковецький

# Зміст

Вступ . . . . .	1
Про стан розвитку відкритих даних в Україні . . . . .	2
Що таке відкриті дані? . . . . .	9
Визначення відкритих даних . . . . .	9
Класифікація відкритих даних . . . . .	9

# Вступ

Мета курсу - ознайомити представників державних структур та представників громадського сектору (активістів, членів громадських організацій, журналістів тощо) з основними поняттями, принципами та підходами щодо роботи з даними з метою їх підготовки та викладення у форматах, що відповідають принципам відкритості та зручності обробки програмними засобами, показати переваги окремих, найбільш популярних, форматів відкритих даних, з огляду на міжнародну практику роботи з даними, а також особливості роботи з відкритими даними в Україні.

Також ці методичні матеріали будуть цікавими представникам ІТ сектору та бізнесу, які мають намір аналізувати публічні дані для працювати з суспільно-корисними ініціативами.

# Про стан розвитку відкритих даних в Україні

За останні два роки Україна досягла великого прогресу в сфері відкритих даних, але все ще знаходиться досить далеко в міжнародних рейтингах.

*Запущено єдиний державний портал відкритих даних [data.gov.ua](http://data.gov.ua)*

Перша версія portalу була розроблена волонтерами із організації Socialboost за підтримки міжнародних організацій і компанії Microsoft. На сьогоднішній день портал переданий на баланс Державного агентства з питань електронного урядування. На порталі доступно майже 6000 наборів даних від 700+ розпорядників, хоча якісних наборів даних у відсотковому співвідношенні не дуже багато.

*Запущено систему державних закупівель ProZorro*

[ProZorro](https://prozorro.gov.ua)<sup>1</sup> – електронна система публічних закупівель яка прийшла на зміну паперовим держтендерам.

У Законі України «Про публічні закупівлі» передбачається:

- Запровадження обов'язковості проведення процедур через електронну систему. (Перший етап — обов'язковість проведення процедур через електронну систему поширюється на головних розпорядників коштів та монополістів (з 1 квітня 2016 року), на другому (з 1 серпня) – на всіх замовників;
- Запровадження електронного аукціону, який передбачає автоматичну оцінку тендерних пропозицій;
- Визначення нових понять «авторизований електронний майданчик», «електронна система закупівель», «централізована закупівельна організація», «система хмарних обчислень»;
- Замість 5 процедур залишити 3 (відкриті торги, конкурентний діалог, переговорна процедура);
- Зміна термінології, зокрема, замість терміну «державна закупівля» вводиться термін «публічна закупівля»; замість термінів «конкурс», «документація конкурсних торгів», «пропозиція конкурсних торгів», «комітет з конкурсних торгів», вводяться поняття «тендер», «тендерна документація», «тендерна пропозиція», «тендерний комітет».

---

<sup>1</sup><https://prozorro.gov.ua>

ProZorro отримала міжнародну премію у сфері публічних закупівель Public Sector Procurement Award за створення і впровадження електронної системи з унікальною архітектурою. Розробка цієї системи на базі програмного забезпечення з відкритим кодом була здійснена у партнерстві влади, бізнесу та громадськості та адмініструвалася антикорупційною організацією [Transparency International Україна](#)<sup>2</sup>. (с)

З 1 серпня 2016 року ProZorro є обов'язковою системою для всіх державних замовників при закупівлі від 200 тис. грн для товарів і від 1,5 млн грн для робіт.

#### *Запущено офіційний портал публічних фінансів України Є-Data*

Є-Data — це офіційний державний інформаційний портал у мережі Інтернет, на якому оприлюднюється інформація про використання публічних коштів та реалізується ідея «Прозорого бюджету». Задача - забезпечити повну прозорість державних фінансів та задовольнити право громадськості на доступ до інформації.

Основні законодавчі та нормативні засади для створення проекту:

- [Закон України «Про відкритість використання публічних коштів»](#)<sup>3</sup>
- План заходів з виконання Програми діяльності Кабінету Міністрів України та Стратегії сталого розвитку “Україна - 2020” у 2015 році п. 95: «Впровадження інтегрованої інформаційно-аналітичної системи “Прозорий бюджет” з метою забезпечення доступності інформації про державні фінанси для суспільних потреб із забезпеченням відкритої звітності за всіма коштами, використаними отримувачами бюджетних коштів..»
- Коаліційна угода Парламенту VIII скликання п. 2.5.10.: «Запровадження системи «Прозорий бюджет» з метою забезпечення доступності інформації про державні фінанси для суспільних потреб.

З 15 вересня 2015 року на порталі оприлюднюються всі трансакції Державної казначейської служби, з листопада на ньому доступна інформація про використання коштів державного і місцевих бюджетів, а у січні інформацію почали розкривати суб'єкти господарювання державної і комунальної власності, у статутному капіталі яких державна або комунальна частка акцій (часток, паїв) перевищує 50 відсотків.

До 2018 року планується забезпечити повний функціонал Інтегрованої інформаційно-аналітичної системи «Прозорий бюджет», яка своєю чергою торкнеться змін бюджетних процесів Міністерства фінансів, автоматизації систем ДФС та Казначейства, автоматизації систем обліку та звітності на місцевих рівнях.

#### *Запущено пошуково-аналітичну систему 007*

<sup>2</sup>[https://uk.wikipedia.org/wiki/%D0%A2%D1%80%D0%B0%D0%BD%D1%81%D0%BF%D0%B5%D1%80%D0%B5%D0%BD%D1%81%D1%96\\_%D0%86%D0%BD%D1%82%D0%B5%D1%80%D0%BD%D0%B5%D1%88%D0%BD%D0%BB](https://uk.wikipedia.org/wiki/%D0%A2%D1%80%D0%B0%D0%BD%D1%81%D0%BF%D0%B5%D1%80%D0%B5%D0%BD%D1%81%D1%96_%D0%86%D0%BD%D1%82%D0%B5%D1%80%D0%BD%D0%B5%D1%88%D0%BD%D0%BB)

<sup>3</sup><http://zakon3.rada.gov.ua/laws/show/183-19>

**Пошуково-аналітична система 007<sup>4</sup>** – це веб-сервіс на основі відкритих даних про використання публічних коштів. Проект передбачає сервіс пошуку та візуалізації даних з відкритих джерел про використання державою бюджетних коштів. Його ідея – дати громадськості інструмент контролю влади та можливість стежити за бюджетними витратами, збирати докази зловживань і швидко переводити боротьбу з корупцією в правове поле. Основний акцент зроблено на простоті використання та представлення специфічної інформації з масивів великих даних. Сайт був відкритий 8 квітня 2016 року.

#### *Нормативно-правове забезпечення*

9 квітня 2015 року Верховна Рада України ухвалила Закон України № 319 «**Про внесення змін до деяких законів України щодо доступу до публічної інформації у формі відкритих даних<sup>5</sup>**». Зазначеним Законом внесені зміни до Закону України «Про доступ до публічної інформації» з метою визначення базових норм та засад розвитку відкритих даних в Україні, а саме:

1. Публічна інформація у формі відкритих даних - це публічна інформація у форматі, що дозволяє її автоматизоване оброблення електронними засобами, вільний та безоплатний доступ до неї, а також її подальше використання;
2. Розпорядники інформації зобов'язані надавати публічну інформацію у формі відкритих даних на запит, оприлюднювати і регулярно оновлювати її на єдиному державному веб-порталі відкритих даних та на своїх веб-сайтах;

3. Будь-яка особа може вільно копіювати, публікувати, поширювати, використовувати, у тому числі в комерційних цілях, у поєднанні з іншою інформацією або шляхом включення до складу власного продукту, публічну інформацію у формі відкритих даних з обов'язковим посиланням на джерело отримання такої інформації.

21 жовтня 2015 року затверджено Постановою КМУ № 835 «**Про затвердження Положення про набори даних, які підлягають оприлюдненню у формі відкритих даних<sup>6</sup>**», якою визначено вимоги до формату і структури наборів даних, а також затверджено перелік пріоритетних наборів даних, які підлягають оприлюдненню (більше 300 наборів). Постановою чітко визначений перелік форматів для оприлюднення відкритих даних в залежності від їх виду:

Текстові дані TXT, RTF, ODT, DOC(X), PDF (с текстовим змістом, скановане зображення), (X)HTML

Структуровані дані RDF, XML, JSON, CSV, XLS(X), ODS, YAML

Графічні дані GIF, TIFF, JPG (JPEG), PNG

Видеодані MPEG, MKV, AVI, FLV, MKS, MK3D

Аудіодані MP3, WAV, MKA

---

<sup>4</sup><https://www.facebook.com/pointOSeven>

<sup>5</sup><http://zakon3.rada.gov.ua/laws/show/319-19>

<sup>6</sup><http://zakon5.rada.gov.ua/laws/show/835-2015-%D0%BF>

## Дані, розроблені з використанням програми Macromedia Flash SWF, FLV

### Архіви даних ZIP, 7z, Gzip, Bzip2

Також, постановою передбачено, що державні реєстри, які постійно оновлюються, мають бути відкриті за допомогою API.

З метою забезпечення ефективного функціонування Єдиного державного веб-порталу відкритих даних та підвищення відкритості і прозорості діяльності органів виконавчої влади та місцевого самоврядування, Державним агентством з питань електронного урядування в лютому 2016 року розроблено проект постанови Кабінету Міністрів України «Деякі питання оприлюднення публічної інформації у формі відкритих даних», які знаходяться зараз на етапі обговорення.

### Оцінка готовності України до розвитку відкритих даних

З метою визначення поточного стану розвитку та готовності України до розвитку відкритих даних, а також планування подальший дій, проведено оцінку готовності України до розвитку відкритих даних за методикою Всесвітнього банку ODRA.

Оцінка проводилась за восьма напрямками:

1. Зобов'язання Уряду;
2. Політичні і правові засади;
3. Інституційна структура, розподіл відповідальності та спроможність урядових структур;
4. Політика та процедури Уряду стосовно обробки даних;
5. Попит на відкриті дані;
6. Залучення громадського сектору та можливості для відкритих даних;
7. Фінансування політики відкритих даних;
8. Національна технологічна інфраструктура та навички.

Якщо коротко, з готовністю громадськості використовувати відкриті дані все добре, більш-менш хороша ситуація з готовністю уряду відкривати дані, технічною інфраструктурою, але все погано в плані фінансування, індустрією розробки продуктів на базі відкритих даних і захистом приватності. Детальний звіт знаходиться за [посиланням](#)<sup>7</sup> (є звіти на українській і англійській мовах).

### Дорожня карта розвитку відкритих даних

З метою забезпечення комплексного розвитку відкритих даних Державним агентством з питань електронного урядування напрацьовано Дорожню карту розвитку відкритих даних в Україні, яка містить 41 завдання по 5 напрямкам:

1. Підвищення доступності та якості відкритих даних;

---

<sup>7</sup><http://dhrp.org.ua/uk/publikatsii/1071-20160227-ua-publication>



2. Розвиток спроможності органів влади щодо публікації відкритих даних;
3. Посилення ролі відкритих даних в реалізації державної політики;
4. Нормативно-правове забезпечення;
5. Розвиток попиту та спроможності цільових аудиторій щодо використання відкритих даних.

Зазначений документ затверджено наказом Мінрегіону від 04.02.2016 року № 19. Повний текст розміщено за [посиланням](#)<sup>8</sup>.

### Хартія відкритих даних

Наразі Україною ініціюється приєднання до міжнародної Хартії відкритих даних.

Проект Постанови КМУ знаходиться за [посиланням](#)<sup>9</sup>.

Розробка Міжнародної хартії була ініційована представниками урядів Канади, Мексики, Великобританії, впливових міжнародних організацій у травні 2015 року під час міжнародної конференції з питань відкритих даних у Канаді.

Головною метою Міжнародної Хартії відкритих даних є покращення та сприяння співпраці та взаємоузгодженості для прийняття та реалізації спільних принципів, стандартів та кращих практик відкритих даних по всьому світу. Цілями Хартії є поширення демократії, боротьба з корупцією та сприяння економічному зростанню по всьому світу. Хартія визначає 6 головних принципів та шляхи розвитку відкритих даних для країни.

### Світові та українські компетенції й рейтинги

Сьогодні найбільш важливими є наступні два рейтинги оцінювання стану розвитку відкритих даних:

1. [Open Data Barometer](#)<sup>10</sup>. За цим рейтингом Україна посідає 62 місце.
2. [Open Data Index](#)<sup>11</sup>. За цим рейтингом Україна посідає 54 місце (із 122).

Україна співпрацює з такими міжнародними організаціями в сфері відкритих даних:

1. [Open Data Institute](#)<sup>12</sup>. Займається розвитком відкритих даних по всьому світу, формуванням стандартів та єдиних підходів, розвитком компетенцій.
2. [Open knowledge foundation](#)<sup>13</sup>. Займається підтримкою громадських інституцій щодо розвитку відкритих даних, а також розвитком відкритої платформи для побудови порталів відкритих даних SKAN.

<sup>8</sup>[https://drive.google.com/file/d/0B1kGsKt9XV\\_QaFZVaTZiT19aRTA/view](https://drive.google.com/file/d/0B1kGsKt9XV_QaFZVaTZiT19aRTA/view).

<sup>9</sup>[https://drive.google.com/file/d/0B1kGsKt9XV\\_QMEFNnklzRXdaSDA/view?usp=sharing](https://drive.google.com/file/d/0B1kGsKt9XV_QMEFNnklzRXdaSDA/view?usp=sharing)

<sup>10</sup>[http://opendatabarometer.org/data-explorer/?\\_year=2015&indicator=ODB&open=UKR](http://opendatabarometer.org/data-explorer/?_year=2015&indicator=ODB&open=UKR)

<sup>11</sup><http://index.okfn.org/place/ukraine>

<sup>12</sup><http://theodi.org>

<sup>13</sup><https://okfn.org>

3. **Open Data for Development**<sup>14</sup>. Займається підтримкою ініціатив з розвитку відкритих даних по всьому світу, а також організацію міжнародної співпраці.

В останньому рейтингу E-Government Development Index (EGDI) 2016 Україна зайняла 62 місце серед 193 країн, покращивши свою позицію на 25 пунктів.

### **Професійні українські організації і ініціативи**

В Україні запущений інкубатор відкритих даних 1991, який системно займається відбором проектів, їх інкубацією і пошуком інвестицій. На сьогодні уже було два набори в інкубаційну програму.

На початку року розпочався EGAP Challenge – конкурс IT проектів в області електронної демократії, соціальної сфери і проектів в сфері відкритих даних. Спільна ініціатива Державного агентства з питань електронного урядування, Фонду Східної Європи в рамках реалізації Програми EGAP, що фінансується Швейцарською Конфедерацією, має на меті не тільки запровадити нові інструменти e-democracy в чотирьох регіонах України (Вінницький, Волинський, Дніпровський та Одеській областях), а також показати українським стартапам нову нішу, в якій можливо створювати якісні проекти, що мають змогу безпосередньо впливати на місцеву і центральну влади. А влада – у свою чергу – отримає набір нових інструментів для росту ефективності її роботи і прозорості взаємодії з платниками податків.

Пріоритетними є напрямки:

- створення нових інструментів взаємодії влади і суспільства, особливо в частині надання можливості громадянам напряму впливати на процеси прийняття управлінських рішень;
- підвищення прозорості і відкритості діяльності органів влади, особливо в частині формування і виконання бюджетів, надання дозвільних документів і тому подібне;
- створення нових якісних сервісів для громадян і бізнесу, особливо в частині надання публічних послуг;
- розвиток проектів на базі відкритих даних;
- вирішення соціальних проблем;
- об'єднання і налагодження ефективного співробітництва громадян для вирішення загальних проблем;
- проекти в області Smart City;
- галузеві проекти, спрямовані на підвищення ефективності державного управління і обслуговування громадян та бізнесу (е-екологія, е-медицина, е-освіта і т.д.).

В чотирьох областях пройшли креативні уїкенди на яких було відібрано 15 проектів для проходження двомісячної інкубації. Партнерами конкурсу виступили компанії Cisco, IBM, DeNovo, Intel.

---

<sup>14</sup><http://od4d.net>

В Україні діють потужні громадські організації, що працюють з відкритими даними— ОПОРА, Чесно, Канцелярська сотня, Vox Ukraine і інші. Результат — десятки досліджень, проекти по візуалізації даних, моніторинг відкритих наборів даних, десятки заходів по усій країні. Декілька проектів уже імплементовані в державні служби і проекти “смарт сіті”.

Наприклад, відкриті дані Укрзалізниці дали можливість провести масштабне дослідження – хто куди подорожує, які ключові станції, як розподіляється потік пасажирів і інші візуалізації.

### Рекомендації по роботі з відкритими даними

На сайті Верховної Ради є [інструкція](#)<sup>15</sup> по роботі з відкритими даними.

Також спільно з Агентством з питань електронного врядування та ПРООН були розроблені [методичні рекомендації щодо оприлюднення наборів даних у формі відкритих даних](#)<sup>16</sup>.

Метою рекомендацій є ознайомлення відповідальних осіб розпорядників інформації з ключовими питаннями, що постають у процесі оприлюднення наборів даних у формі відкритих даних. Також, методичні рекомендації будуть корисні розробникам застосунків на базі відкритих даних, громадським організаціям, засобам масової інформації та громадянам, які мають намір працювати з відкритими даними.

Проект першої версії Методичних рекомендацій підготовлено на основі найбільш популярних загальних запитань від розпорядників інформації, але, безумовно, ще не охоплює усіх піднятих розпорядниками інформації та громадськістю питань. Методичні рекомендації будуть систематично оновлюватися з метою забезпечення розпорядників інформації необхідною інформацією для оприлюднення наборів даних.

---

В цілому, Україна знаходиться лише на початку з точки зору відкритості державних служб, створення інститутів моніторингу діяльності депутатів, чиновників і міських служб, дигіталізації державних послуг. За наявності політичної волі, допомоги міжнародних організацій, а також за умови координації громадських організацій, бізнесу та ІТ сектору Україна здатна на досить короткий термін значно поліпшити свої показники в області електронного врядування та відкритих даних.

---

<sup>15</sup><http://data.rada.gov.ua/open/main/opendata>

<sup>16</sup>[https://drive.google.com/file/d/0B1kGsKt9XV\\_QWjhaZ0ZMVmFiUE0/view](https://drive.google.com/file/d/0B1kGsKt9XV_QWjhaZ0ZMVmFiUE0/view)

# Що таке відкриті дані?

## Визначення відкритих даних

Відкриті дані (англ. Open Data) – це концепція, яка відображує ідею, що визначені дані мають бути доступні для легкої обробки програмними засобами (machine readable) та подальшого використання і розповсюдження без жодних обмежень і контролю, в тому числі й для комерційного використання. Відкриті дані – це не просто інформація, а концепція, тобто система поглядів, підходів, процесів, які мають одну ідею та мету – вільного використання і розповсюдження даних про діяльність державних органів та органів місцевого самоврядування через мережу Інтернет.

Згідно загальноприйнятої в світі концепції, обов'язково безоплатними можуть бути тільки ті дані, які знаходяться у власності держави («ліцензійно чисті»), і якщо вони подаються в первинному необробленому вигляді. Додаткова обробка або доступ до API частіше за все лімітуються або коштують грошей. Крім того, хто завгодно може (з обов'язковим посиланням на джерело) використати дані в комерційний спосіб, створити на базі них власну програму чи обробити та надати нову цінність (розкласифікувати, встановити зв'язки тощо). Безумовно, всі витрати на додаткову роботу редакторів, програмістів та інше оплачує замовник (або кінцевий користувач).

Відкриті дані дозволяють повторно і необмежено використовувати інформацію, поєднувати її між собою, зменшити або виключити зайві витрати на дублювання та опрацювання великих масивів даних, реєстрів, довідників, баз даних тощо, створених в різних органах влади.

Закон “Про внесення змін до деяких законів України щодо доступу до публічної інформації у формі відкритих даних” №2171 від 19.02.2015 (прийнятий 09.04.2015) визначає наступний термін:

Публічна інформація у формі відкритих даних – це публічна інформація у форматі, що дозволяє її автоматизоване оброблення електронними засобами, вільний та безоплатний доступ до неї, а також її подальше використання.

## Класифікація відкритих даних

В законопроекті одночасно використовуються поняття «форма» та «формат», що не є тотожними. Розглянемо детально.

Для того, щоб зрозуміти, які можуть бути форми відкритих даних, ми також звернемося до відомої класифікації «5 зірок Open Data» (<http://5stardata.info/><sup>17</sup>), де якість даних та рівень відкритості визначається кількістю зірок від 1 до 5, чим більше – тим краще. Відкритість даних залежить від способів доступу, форматів та кількості додаткових дій, які потрібні для отримання кінцевої інформації, її обробки та збереження у власному сховищі або базі даних.

## 5 Star Data Schema

★ Available on the web (whatever format) but with an open license, to be Open Data

★★ Available as machine-readable structured data (e.g. excel instead of image scan of a table)

★★★ as (2) plus non-proprietary format (e.g. CSV instead of excel)

★★★★ All the above plus, Use open standards from W3C (RDF and SPARQL) to identify things, so that people can point at your stuff

★★★★★ All the above, plus: Link your data to other people's data to provide context



5-Star Open Data Scheme  
September 19, 2014



Одну зірку (\*) отримує будь-яка інформація вільно доступна через Інтернет в будь-якому форматі. Під цю класифікацію підпадає файл в форматі PDF або інша (сканована) копія документу, на який веде пряме посилання на офіційному сайті державного органу. Якщо цей файл можна відкрити на власному екрані, прочитати, роздрукувати та отримати звітні потрібну інформацію, то це відкриті дані з однією зіркою.

Дві зірки (\*\*) отримує структурована інформація, яку можна обробляти автоматично, наприклад, в форматах для веб-браузерів чи офісних програм (відкриті формати – TXT, HTML, RSS; пропріетарні формати, Excel – XLS, Word – DOC, RTF). Якщо дані знаходяться в тілі вихідної веб-сторінки, але не мають чіткої структури, містять зайві елементи оформлення, навігації, якщо дані потребують додаткових дій – спеціального розбору (парсингу), то вони вважаються «з двома зірками».

---

<sup>17</sup><http://5stardata.info/>

**Три зірки (\*\*\*)** може отримати інформація, представлена у відомих, добре описаних відкритих структурованих форматах (наприклад, CSV, JSON, XML, YAML) і якщо автоматизована її обробка не потребує від користувача особливих ліцензій та додаткових плат. До відкритих форматів також відносяться пов'язані дані (HTML+RDFa) з узгодженою розміткою елементів в атрибутах (див. сніппет для пошукових систем) або текстові файли таблиць, поля яких розділені табуляцією, комами, крапками з комою або іншими символами.

**Чотири зірки (4)** надаються у випадку, якщо можна отримати первинні необроблені набори відкритих даних у вигляді файлів (довідники, списки, таблиці у відкритому форматі, зліпок бази даних, архів документів тощо) або фільтровані дані у запиті до API за вказаними параметрами. Це дає змогу отримувати тільки потрібну інформацію, актуальну на момент запиту, заощаджує ресурси та час користувача. Безумовно, API має бути описаний так само, як і формати даних, а доступ до нього може бути анонімний без обмежень або з реєстрацією, за вказаним ідентифікатором, лімітами на кількість одночасних запитів тощо.

Останній рівень – **п'ять зірок (5)** – надається інформації, коли набори відкритих даних пов'язані між собою (мають спільні довідники, класифікатори, ідентифікатори, посилання між документами та іншими елементами тощо) і представляють собою семантичну мережу, що постійно оновлюється й змінюється відповідно до сучасних запитів.

## Формати даних

В залежності від специфіки даних, їх розміру та тематики (геологія, тендери, реєстри, судові документи тощо), одні проекти відкритих даних створювались на базі наборів PDF чи DOC файлів, таблиць XLS, що перетворювались на прості текстові таблиці CSV, а інші брали за основу формат розмітки XML, проектували власні схеми XSD і використовували складні структури.

Як свідчить остання статистика використання форматів відкритих даних, найбільш поширений в світі формат (як по кількості, так і по об'єму даних) – PDF. Для українських органів влади, де найбільш розповсюджені операційні системи Microsoft Windows, переважають формати DOC та XLS. Разом з новими версіями офісних програм в Інтернет почали з'являтися документи DOCX та XLSX, рідко ODF (*Open Document Format*). Як надбання DOS'івського періоду, поки ще зустрічаються документи в старому форматі Lexicon (TXT) або ранніх версій Word (RTF). Після поширення ініціативи відкриття державних даних та створення порталів, кількість наборів в форматі XML та інших відкритих форматах почала суттєво збільшуватись.

Необроблені дані, сформовані державними структурами за багато років, можуть бути досить неоднорідними, а деякі набори навіть дублюються в різних форматах для зручності користування.

Серед доступних в цей час форматів відкритих даних, які можна автоматично обробляти електронними засобами, є: CSV (текстові дані, розділені комами або іншими розділовими

символами), JSON (формат, орієнтований на обробку складних структурованих даних за допомогою javascript) та XML (універсальний текстовий формат розмітки). Є багато форматів, що основані на XML, зокрема, HTML, KML та інші.

## Табличні формати

### CSV (Comma-Separated Values)

CSV (від англ. *Comma-Separated Values* – значення, що розділені комами) – текстовий відкритий формат, призначений для представлення таблиць (масивів, наборів) даних, де кожний рядок – це запис таблиці, а значення окремих полів у рядку розділені спеціальними символами (delimiter), зазвичай комами. Щоб завантажити записи таблиці за найменуванням полів (а не за порядковими номерами), додатково потрібно мати опис її структури – назви та формат полів.

Більшість програм широко трактують цей формат і допускають використання інших розділових символів, наприклад, табуляції (TSV) чи коми з крапкою.

Приклад даних в форматі CSV:

- 1 1997,Ford,E350,"ac, abs, moon",3000.00
- 2 1999,Chevy,"Venture ""Extended Edition""",",4900.00
- 3 1996,Jeep,Grand Cherokee,"MUST SELL! air, moon roof, loaded",4799.00

Таблиця в форматі CSV має наступні параметри: кодування (Windows-1251, KOI-8, UTF-8 тощо) і формат кінця рядків. Різниця у символах розділення рядків залежить від операційної системи (для Unix – один символ CR (CHR 0x0D), для Mac – LF (CHR 0x0A), а для Windows – пара символів CR LF).

Інститут відкритих даних Open Data Institute випустив оновлену версію 0.4 додатка Comma Chameleon <https://github.com/theodi/comma-chameleon/releases/tag/0.4.0><sup>18</sup> - з валідації CSV файлів. Ця версія найбільш стабільна з усіх попередніх і існує в версіях для Mac, Linux, Windows і просто як відкритий код в репозиторії <https://github.com/theodi/comma-chameleon><sup>19</sup>.

Корисний інструмент для всіх хто готує дані для публікації і думає про автоматизацію і спрощення очищення даних.

Корисні сервіси та інструменти для роботи з CSV файлами:

1. Валідація CSV файлів: Comma Chameleon <https://github.com/theodi/comma-chameleon><sup>20</sup>.
2. CSVLint <http://csvlint.io/><sup>21</sup> - онлайн сервіс з валідації CSV файлів і з відкритим кодом <https://github.com/theodi/csvlint><sup>22</sup>.

<sup>18</sup><https://github.com/theodi/comma-chameleon/releases/tag/0.4.0>

<sup>19</sup><https://github.com/theodi/comma-chameleon>

<sup>20</sup><https://github.com/theodi/comma-chameleon>

<sup>21</sup><http://csvlint.io/>

<sup>22</sup><https://github.com/theodi/csvlint>

3. CSVkit <https://github.com/wireservice/csvkit><sup>23</sup> - бібліотека для Python за численними маніпуляціям з CSV файлами і безліччю інструментів для командного рядка.
4. Textql <https://github.com/dinedal/textql><sup>24</sup> - інструмент із запуску SQL запитів на CSV / TSV файлах.
5. PapaParse <https://github.com/mholt/PapaParse><sup>25</sup> - парсер дуже великих CSV файлів.
6. <https://mledoze.github.io/countries/><sup>26</sup> - країни світу в JSON, CSV, XML і YAML.
7. Tablib <https://github.com/kennethreitz/tablib><sup>27</sup> - бібліотека для роботи з будь-якими табличними даними включаючи CSV.

## XLS/XLSX (Document Office Open XML)

Файл XLSX - електронна таблиця, створена в Microsoft Excel - додатку для роботи з таблицями. Дані в документі зберігаються в комірках, кожна з яких має певну адресу (колонки нумеруються англійськими літерами, рядки - цифрами, наприклад: A1 - ліва верхня клітинка).

Кожна клітинка може містити як фіксовані дані, так і формули, часто пов'язані з даними в інших осередках. Формат XLSX дозволяє користувачеві змінювати форматування тексту: його шрифт, колір, вирівнювання (в комірці) і інші параметри. Крім того, документ може містити зображення, а також діаграми, побудовані на основі даних в певних комірках.

Програми, які дозволяють працювати з XLSX файлами: File Viewer Lite, Microsoft Excel, Microsoft Excel Viewer, Microsoft Works, Corel WordPerfect Office X7, Apache OpenOffice, Kingsoft Spreadsheets, Nuance OmniPage Ultimate, Gnumeric, LibreOffice.

Спочатку формат XLSX створювався як заміна раніше бінарного формату документів (.xls), який раніше використовували додатки Microsoft Office. Але в 2006 році формат Office Open XML був оголошений вільним і відкритим форматом.

## Графічні формати

### JPEG (Joint Photographic Experts Group)

JPEG (англ. *Joint Photographic Experts Group*, за назвою організації-розробника) - один з популярних графічних форматів, застосовуваний для зберігання фотозображень і подібних до них зображень. Файли даних JPEG, зазвичай мають розширення .jpg, .jfif, .jpe або .jpeg. Однак з них .jpg є найпопулярнішим на всіх платформах.

Алгоритм JPEG дозволяє стискати зображення як з втратами, так і без втрат (режим стиснення lossless JPEG). Підтримуються зображення з лінійним розміром не більше 65535 × 65535 пікселів.

---

<sup>23</sup><https://github.com/wireservice/csvkit>

<sup>24</sup><https://github.com/dinedal/textql>

<sup>25</sup><https://github.com/mholt/PapaParse>

<sup>26</sup><https://mledoze.github.io/countries/>

<sup>27</sup><https://github.com/kennethreitz/tablib>



Алгоритм JPEG найбільшою мірою придатний для стиснення фотографій і картин, що містять реалістичні сцени з плавними переходами яскравості і кольору. Найбільшого поширення JPEG отримав в цифровій фотографії і для зберігання і передачі зображень з використанням мережі Інтернет.

Формат JPEG в режимі стиснення з втратами малоприматний для стиснення креслень, текстової та знакової графіки, де різкий контраст між сусідніми пікселями приводить до появи помітних артефактів. Такі зображення доцільно зберігати в форматах без втрат, таких як JPEG-LS, TIFF, GIF, PNG або використовувати режим стиснення Lossless JPEG.

JPEG не підходить для стиснення зображень багатоетапної обробки, так як спотворення в зображення будуть вноситися кожен раз при збереженні проміжних результатів обробки. JPEG не повинен використовуватися і в тих випадках, коли недопустимі навіть мінімальні втрати, наприклад, при стисненні астрономічних або медичних зображень.

Детальніше про формат: <https://ru.wikipedia.org/wiki/JPEG><sup>28</sup>

## PNG (portable network graphics)

PNG (англ. *Portable network graphics*) - растровий формат зберігання графічної інформації, що використовує стиснення без втрат за алгоритмом Deflate. Формат PNG розроблений для заміни форматів GIF та TIFF. Формат PNG позиціонується передусім для використання в Інтернеті і редагування графіки. Формат PNG зберігає графічну інформацію в стислому вигляді, причому це стиснення графічної інформації проводиться без втрат.

PNG підтримує три основних типи растрових зображень:

- Півтонування (з глибиною кольору 16 біт)
- Кольорові індексовані зображення (палітра 8 біт для кольору глибиною 24 біт)
- Повнокольорове зображення (з глибиною кольору 48 біт)

Формат має наступні основні переваги:

- практично необмежену кількість кольорів в зображенні;
- опціональна підтримка альфа-каналу;
- можливість гамма-корекції;
- двовимірна черезстрокова розгортка;
- можливість розширення формату для користувача блоками (на цьому заснований, зокрема, APNG).
- будь-яке збережене зображення PNG може бути прочитано в будь-якому іншому додатку, що підтримує PNG.

---

<sup>28</sup><https://ru.wikipedia.org/wiki/JPEG>

PNG є хорошим форматом для редагування зображень, навіть для зберігання проміжних стадій редагування, так як відновлення і повторне збереження зображення здійснюється без втрат в якості.

Детальніше про формат: <https://ru.wikipedia.org/wiki/PNG><sup>29</sup>

## TIFF (Tagged Image File Format)

TIFF (англ. *Tagged Image File Format*) - формат зберігання растрових графічних зображень. TIFF став популярним форматом для зберігання зображень з великою глибиною кольору. Він використовується при скануванні, відправленні факсів, розпізнаванні тексту, в поліграфії, широко підтримується графічними додатками. TIFF був обраний в якості основного графічного формату операційної системи NeXTSTEP і з неї підтримка цього формату перейшла в Mac OS X.

Спочатку формат підтримував стиснення без втрат, згодом формат був доповнений для підтримки стиснення з втратами в форматі JPEG. Файли формату TIFF, як правило, мають розширення .tiff або .tif.

Структура формату гнучка і дозволяє зберігати зображення в режимі кольорів з палітрою, а також в різних колірних просторах:

- Бінарному (двокольоровому, іноді неправильно званому чорно-білим)
- напівтоновому
- з індексованої палітрою
- RGB
- CMYK
- YCbCr
- CIE Lab

Підтримуються режими 8, 16, 32 і 64 біт на канал при цілочисельному, а також 32 і 64 біт на канал при поданні значення пікселя числами з плаваючою комою.

TIFF є тяговим форматом і в ньому є такі види міток:

- основні мітки;
- розширені мітки;
- спеціальні мітки.

Детальніше про формат: <https://ru.wikipedia.org/wiki/TIFF><sup>30</sup>

---

<sup>29</sup><https://ru.wikipedia.org/wiki/PNG>

<sup>30</sup><https://ru.wikipedia.org/wiki/TIFF>

## Текстові формати

### TXT (Textfile)

TXT - це формат, що містить текстові дані, які, як правило, організовані в вигляді рядків; служить основою для багатьох більш спеціалізованих форматів, таких як CHM, PHP, XML, CSV. Може бути переглянутий будь-яким текстовим редактором.

Текстові файли розбиваються на кілька рядків. На сучасних платформах розбивка на рядки кодується керуючим символом переведення рядка, а іноді послідовністю двох керуючих символів. Взагалі, текстові файли можуть містити друковані символи, такі як букви, цифри і розділові знаки і деяку кількість керуючих символів, таких як знаки табуляції і переведення рядка.

Файли TXT містять дуже мало елементів форматування, однак дозволяють співвіднести прийняті набори елементів форматування з системним терміналом або простим текстовим редактором. Файли TXT універсальні за своєю природою, тому що відкрити такі файли в стані будь-який текстовий редактор. При цьому файли TXT можуть використовувати Unicode, яка дозволяє полегшити користувачам використання файлів TXT, які пишуть на різних мовах. Файли TXT, що використовують текст тільки в кодуванні ASCII, можна переносити на різні комп'ютери і пристрої на ОС Unix, Mac і Windows.

### Markdown

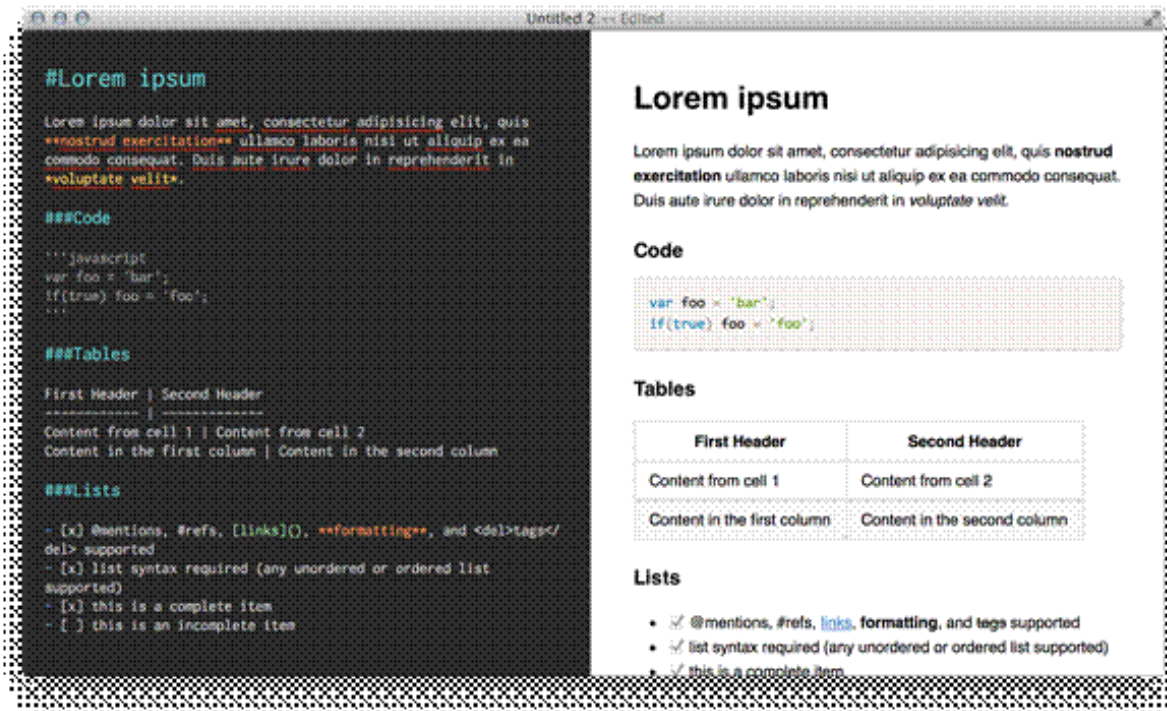
Markdown - полегшена мова розмітки, створена з метою написання максимально читабельного і зручного для редагування тексту, але придатного для перетворення в мови для публікацій (HTML, Rich Text і ін.).

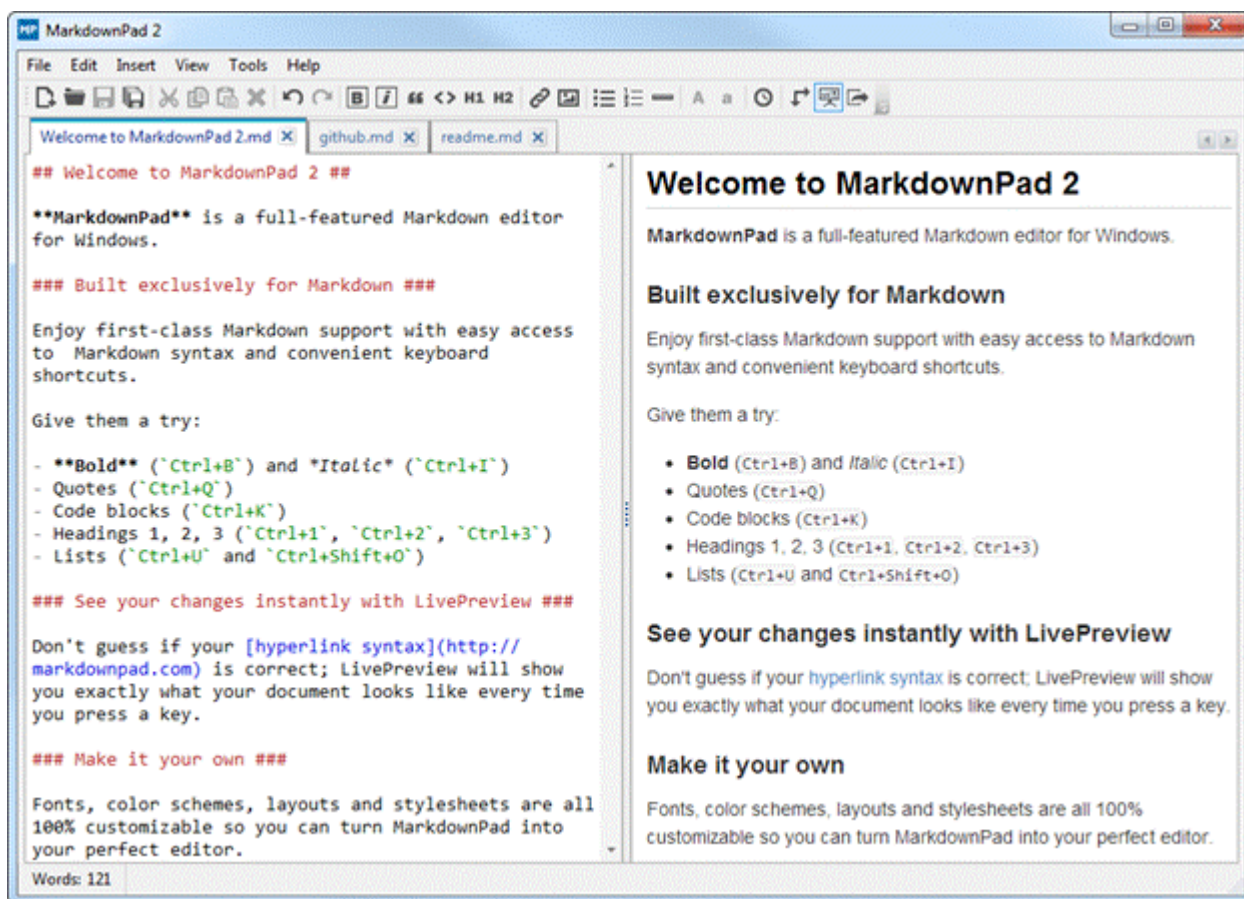
Таким чином, "Markdown" це дві речі:

- простий синтаксис форматування тексту;
- програмний інструмент, написаний на Perl, який перетворює звичайний форматування тексту в HTML.

Першочерговою метою розробки для форматування синтаксису Markdown є зробити його максимально читабельним, наскільки це можливо. Ідея полягає в тому, що Markdown-форматований документ повинен бути опублікований як звичайний текст, використовуючи теги та інструкції форматування.

Приклади форматування тексту мовою розмітки Markdown:





Детальніше про формат та його синтаксис:

1. Markdown <http://daringfireball.net/projects/markdown/><sup>31</sup>
2. Wikipedia <https://uk.wikipedia.org/wiki/Markdown><sup>32</sup>

## Текстово-графічні формати

### HTML (HyperText Markup Language)

HTML (від англ. *HyperText Markup Language* - «мова гіпертекстової розмітки») - стандартна мова розмітки документів у Всесвітній павутині. Більшість веб-сторінок містять опис розмітки на мові HTML (або XHTML). Мова HTML інтерпретується браузерами; отриманий в результаті інтерпретації форматований текст відображається на екрані монітора комп'ютера або мобільного пристрою.

У всесвітній павутині HTML-сторінки, як правило, передаються браузерам від сервера по протоколам HTTP або HTTPS, у вигляді простого тексту або з використанням шифрування.

<sup>31</sup><http://daringfireball.net/projects/markdown/>

<sup>32</sup><https://uk.wikipedia.org/wiki/Markdown>

Текстові документи, що містять розмітку на мові HTML (такі документи зазвичай мають розширення .html або .htm), обробляються спеціальними додатками, які відображають документ в його форматованому вигляді. Такі додатки, звані «браузерами», зазвичай надають користувачеві зручний інтерфейс для запиту веб-сторінок, їх перегляду (і виведення на інші зовнішні пристрої) і, при необхідності, відправки введених користувачем даних на сервер.

HTML - тегова мова розмітки документів. Будь-який документ на мові HTML являє собою набір елементів, причому початок і кінець кожного елемента позначається спеціальними позначками - тегами. Елементи можуть бути порожніми, тобто не містять ніякого тексту та інших даних (наприклад, тег переносу рядка <br>). В цьому випадку зазвичай не вказується закриваючий тег. Крім того, елементи можуть мати атрибути, що визначають будь-які їх властивості (наприклад, розмір шрифту для елемента font). Атрибути вказуються в відкриваючому тегові.

Детальніше про формат: <https://ru.wikipedia.org/wiki/HTML><sup>33</sup>

## DOCX (Document Office Open XML)

DOCX (*Document Office Open XML*) - формат файлу для зберігання електронних документів пакетів офісних додатків - зокрема, Microsoft Office. Формат є zip-архів, що містить текст у вигляді XML, графіку і інших даних, які можуть бути переведені в послідовність бітів (серіалізовані) із застосуванням захищених патентами довічних форматів, специфікації яких були опубліковані Microsoft для користувачів на умовах Microsoft Open Специфікація Promise (англ.).

Спочатку формат створювався як заміна раніше бінарного формату документів, який використовували додатки Microsoft Office аж до версії Office 2003 включно. У 2006 році формат Office Open XML був оголошений вільним і відкритим форматом Ecma International. Він є форматом за умовчанням для додатків Microsoft Office 2007 і пізніших.

Файл з розширенням .DOCX - документ, створений Microsoft Word, програмою обробки тексту; містить текст документа, зображення, форматування, стилі, намальовані об'єкти та інші параметри документа.

Файли DOCX створюються за допомогою відкритого формату XML, в якому зберігаються документи, як збори окремих файлів і папок в стислому пакеті. DOCX-файли містять XML-файли і три папки, docProps, Word, і \_rels, які містять властивості документа, зміст і відносини між файлами.

Детальніше про формат: [https://ru.wikipedia.org/wiki/Office\\_Open\\_XML](https://ru.wikipedia.org/wiki/Office_Open_XML)<sup>34</sup>

## OpenDocument format (LibreOffice/OpenOffice)

Формат OpenDocument - це формат для текстових, табличних документів та презентацій, що був створений з метою підвищити інтероперабельність між різними офісними застосунками

---

<sup>33</sup><https://ru.wikipedia.org/wiki/HTML>

<sup>34</sup>[https://ru.wikipedia.org/wiki/Office\\_Open\\_XML](https://ru.wikipedia.org/wiki/Office_Open_XML)

(програмами, офісними пакетами).

Історична довідка: процес стандартизації форматів офісних документів почався на початку 2000-х років та був обумовлений становищем на ринку офісних пакетів програм. Microsoft на той момент займала домінантне становище. Європейська комісія за скаргою Sun Microsystems, що тоді розробляла альтернативний офісний пакет OpenOffice.org, вирішила, що Microsoft зловживає своїм становищем та веде нечесну конкуренцію. Microsoft змушена була відкрити специфікацію своїх форматів. Паралельно з цим, відбувалась стандартизація офісних форматів. Результатом зусиль зі стандартизації стало створення двох конкуруючих форматів - Office Open XML (Microsoft) та OpenDocument (Sun Microsystems).

Файли формату OpenDocument мають розширення ODT (OpenDocument Text), ODS (OpenDocument Spreadsheet). Його підтримують наступні пакети програм:

- Apache OpenOffice
- LibreOffice
- MS Office (2013 - тільки читання, 2016 - читання та редагування)
- Google Docs (тільки імпорт/експорт)

## PDF (Portable Document Format)

Формат переносного документа (PDF) - це формат файлу, який використовується для надійного уявлення і обміну документами, незалежно від програмного і апаратного забезпечення або операційної системи. В даний час формат PDF, винайдений компанією Adobe, є відкритим стандартом, підтримуваним Міжнародною організацією зі стандартизації (ISO). Файли PDF можуть містити посилання і кнопки, поля форм, аудіо- та відеоелементи, а також бізнеслогіку. Вони також підтримують можливість електронного підпису і можуть переглядатися за допомогою безкоштовного програмного забезпечення Acrobat Reader DC.

Можливості формату:

- мають такий же вигляд, як вихідні документи, незалежно від використовуваної платформи або пристрою;
- відповідають стандартам ISO 32000 з обміну електронними документами, включаючи спеціальні стандарти для PDF;
- можуть бути підписані електронно будь-яким користувачем за допомогою безкоштовного програмного забезпечення Adobe Acrobat Reader DC або мобільного застосування Acrobat DC;
- можуть захищатися паролем для запобігання копіювання та редагування;
- можуть виправлятися для видалення конфіденційної інформації без можливості відновлення;
- зберігають всю інформацію вихідного файлу, включаючи текст, зображення, аудіо, 3D-карти і багато іншого, об'єднану в одному файлі;



- забезпечують зручний пошук, включаючи пошук по сканованому тексту, який був перетворений за допомогою технології оптичного розпізнавання Символів (OCR);
- працюють з допоміжними технологіями, які роблять файли PDF доступнішими для людей з обмеженими можливостями, такими як слабкий зір або сліпота.

Детальніше про формат: [https://ru.wikipedia.org/wiki/Portable\\_Document\\_Format](https://ru.wikipedia.org/wiki/Portable_Document_Format)<sup>35</sup>  
та <https://acrobat.adobe.com/ru/ru/products/about-adobe-pdf.html><sup>36</sup>

## Формати представлення даних через API

### JSON (JavaScript Object Notation)

JSON (від англ. *JavaScript Object Notation*) – текстовий відкритий формат, оснований на Javascript представлені та призначений для обміну даними в Інтернет між сервером та клієнтом або сервером і сервером. Хоча він позиціонується, як незалежний від системи і мови програмування, частіше за все використовується за допомогою програм на Javascript, але як і інші текстові формати, легко читається людиною.

Найбільшу популярність JSON набув після створення інтерактивних веб-сторінок, дані до яких через API передавались під час взаємодії користувача з елементами інтерфейсу (т.з. технологія AJAX). За рахунок своєї лаконічності, на відміну від XML, простоті й швидкості використання саме в програмах на Javascript, широкими можливостями в серіалізації даних – рекурсивного перетворення в текстовий вигляд складних об'єктів, формат активно використовується для формування «на льоту» та передачі структур даних в Інтернет в різних інформаційних системах і сервісах.

Синтаксис JSON (див. <https://ru.wikipedia.org/wiki/JSON><sup>37</sup>) представляє собою текст у вигляді програмного коду Javascript, який описує кілька варіантів структур:

- 1) масив в квадратних дужках (значення через кому),
- 2) іменованний масив (чи об'єкт) у фігурних дужках (пари ключ-значення, розділені двома крапками, через кому),
- 3) просте значення (число, рядок, true, false, null),
- 4) функцію.

Для перетворення формату JSON в об'єкт в пам'яті (десеріалізації) в Javascript використовується функція `eval()`, яка виконує завантажений в текстовий рядок програмний код.

Якщо для обробки використовувати тільки мову Javascript, виникають як мінімум дві проблеми. По-перше, щоб обробити дані, потрібно мати приблизно в два рази більше вільної пам'яті, ніж файл (спочатку завантажуються повний текст JSON, а після його виконання, дані об'єкту завантажуються повторно у вигляді відповідних структур). По-друге, можливість

<sup>35</sup>[https://ru.wikipedia.org/wiki/Portable\\_Document\\_Format](https://ru.wikipedia.org/wiki/Portable_Document_Format)

<sup>36</sup><https://acrobat.adobe.com/ru/ru/products/about-adobe-pdf.html>

<sup>37</sup><https://ru.wikipedia.org/wiki/JSON>



передачі програмного коду, який автоматично стартує, дозволяє зловмисникам отримати несанкціонований доступ до системи. Цю вразливість також можна уникнути, якщо спочатку перевірити і очистити текст JSON від заборонених об'єктів, а дані перетворити на коректні значення.

Через нераціональне використання пам'яті і ресурсів в JSON краще працювати зі об'єктами будь-якою складністю, але не дуже великих розмірів, в швидкому режимі «запит-відповідь».

Приклад формату JSON:

```
1 {  
2   "firstName": "Иван",  
3   "lastName": "Иванов",  
4   "address": {  
5     "streetAddress": "Московское ш., 101, кв.101",  
6     "city": "Ленинград",  
7     "postalCode": 101101  
8   },  
9   "phoneNumbers": [  
10    "812 123-1234",  
11    "916 123-4567"]  
12 }
```

## XML (eXtensible Markup Language)

XML (від англ. *eXtensible Markup Language* – мова розмітки, що розширюється) – мабуть найстаріший текстовий відкритий формат, створений в 1994 році та рекомендований Консорціумом Всесвітньої павутини (W3C), як основний для обміну інформацією в Інтернет. Гіпертекстова розмітка (HTML) – це один з різновидів XML. Разом з таблицями каскадних стилів CSS, які формують зовнішній вигляд документів, вони є тими основними форматами, що обумовлюють розвиток технологій.

Насправді XML (див. <https://ru.wikipedia.org/wiki/XML><sup>38</sup>) – універсальний зручний для людини формат для збереження, передачі структурованих даних і їх автоматичної обробки у формі машиночитних документів. З нього починали розвиток інші відкриті формати, на ньому вдосконалювались підходи обміну даними. Сама ж мова розмітки XML є похідною від ще старішого та більш складного формату SGML (Standard Generalized Markup Language – стандартна узагальнена мова розмітки), стандартизованого за ISO 8879:1986 Information processing—Text and office systems—Standard Generalized Markup Language (SGML).

Перевагами XML є простота та гнучкість розмітки, яка не вимагає формальних, фіксованих назв тегів чи параметрів, і будь-який розробник може доповнювати та змінювати формат, створювати власну схему XSD (XML Schema Definition). Фактично, це мова, яка описує

---

<sup>38</sup><https://ru.wikipedia.org/wiki/XML>

сама себе і будь-які за розміром і складністю структури даних. Безумовно, в цей формат можна конвертувати інші формати (наприклад, XLS). Крім того, формат не залежить від операційної системи чи бази даних. Для простого перегляду чи редагування достатньо текстового редактора. Але є й недоліки – дані в форматі XML займають значно більше місце, ніж це потрібно, за рахунок повторення тегів та відступів, а парсинг значно складніший, ніж CSV чи JSON.

За довгий час існування XML на його базі було розроблено багато форматів і стандартів зі схожим синтаксисом (див. приклади стандартних схем на <http://schema.org><sup>39</sup>). Зазвичай цю групу форматів називають загальною назвою – XML, тому що вони мають єдині механізми опису схем XSD, перевірки правильності даних (валідації), перегляд DOM (Document Object Model), доступу до елементів XPath та трансформації для автоматичного конвертування у інші схеми чи формати (наприклад, альтернативні JSON та YAML) за допомогою мови перетворення XSLT (*eXtensible Stylesheet Language Transformations*).

Використання формату XML (а саме LegalXML) у якості відкритого стандарту нормативно-правового документа – це сучасний спосіб забезпечити обмін інформацією (документами, картками, довідниками тощо) між інформаційними системами або в межах однієї системи при підготовці (опрацюванні) документів.

В одному файлі XML в текстовому вигляді, крім основних даних та тексту електронного нормативного документа, можна розміщати метадані (характеристики, реквізити, опис, класифікацію тощо), вкладені файли (картинки, стилі тексту, таблиці Excel, документи Word і т.і.), необхідні структури чи довідники. Це дозволяє зручно не тільки зберігати, передавати, обробляти документ, отримувати PDF версію для друку, формувати зміст чи робити посилання на конкретну главу, статтю, пункт, підпункт тощо, а й автоматизовано вносити зміни, підготовлені у вигляді, що дозволяє їх програмну обробку.

Приклад даних у форматі XML:

```
1 <person firstName="Иван" lastName="Иванов">
2     <address streetAddress="Московское ш., 101, кв.101"
3         city="Ленинград" postalCode="101101"/>
4     <phoneNumbers>
5         <phoneNumber>812 123-1234</phoneNumber>
6         <phoneNumber>916 123-4567</phoneNumber>
7     </phoneNumbers>
8 </person>
```

## RDF (Resource Description Framework)

RDF - це розроблена консорціумом Всесвітньої павутини модель для представлення даних, особливо - метаданих. RDF представляє твердження про ресурсах у вигляді, придатному для машинної обробки. RDF є частиною концепції семантичної павутини.

---

<sup>39</sup><http://schema.org>

Ресурсом в RDF може бути будь-яка сутність - як інформаційна (наприклад, веб-сайт або зображення), так і неінформаційна (наприклад, людина, місто або якесь абстрактне поняття). Твердження, що висловлюється про ресурс, має вигляд «суб'єкт - предикат - об'єкт» і називається триплетом. Затвердження «небо блакитного кольору» в RDF-термінології можна представити таким чином: суб'єкт - «небо», предикат - «має колір», об'єкт - «блакитний». Для позначення суб'єктів, відносин і об'єктів в RDF використовуються URI.

RDF сам по собі є не форматом файлу, а тільки лише абстрактною моделлю даних, тобто описує пропоновану структуру, способи обробки та інтерпретації даних. Для зберігання і передачі інформації, покладеної в модель RDF, існує цілий ряд форматів запису.

Для обробки RDF-даних пропонується реалізувати мови запитів: SPARQL (стандарт W3C), RQL, RDQL.

Для запису і передачі RDF використовується кілька форматів, в тому числі:

- RDF / XML - запис у вигляді XML-документа;
- RDF / JSON - запис у вигляді JSON-даних;
- RDFa - запис всередині атрибутів довільного HTML- або XHTML-документа;
- N-Triples, Turtle, N3 - компактні форми запису тверджень.

## Формати даних для роботи з геопросторовими даними

### Формати GIS

Файли геопросторових даних можуть містити таку інформацію:

- Просторова або геометрична інформація, яка забезпечує положення і форми конкретних географічних об'єктів.
- Інформація про атрибути, яка забезпечує описову інформацію (числову, текстову, логічну) про кожний об'єкт.
- Інформація про відображення, яке може описувати спосіб оформлення об'єкту.

Деякі цифрові карти не містять всі три типи інформації. Наприклад, растрові карти, як правило, не включають в себе інформацію про атрибути. В свою чергу багато джерел векторних даних не включають інформацію про оформлення.

### Просторова або геометрична інформація

Географічна інформація описує у цифровій формі положення та форму кожного зображуваного об'єкта. На векторній карті, властива об'єкту позиція зазвичай виражається у вигляді наборів координат  $x$ ,  $y$  або  $x$ ,  $y$ ,  $z$  з використанням системи координат, визначеній для набору геоданих. Більшість векторних географічних інформаційних систем підтримують три основних геометричні об'єкти: \* Точка: одна пара координат, 0-вимірний об'єкт. \*

Лінія: дві або більше точок у певній послідовності, одновимірний об'єкт. \* Полігон: площа, обмежена лінією, двовимірний об'єкт. \* Тіло, об'єм обмежений гранями, тривимірний об'єкт.

Векторні дані характеризуються топологією, яка задає правила просторової взаємодії між векторними об'єктами. В залежності від масштабу даних один і той же об'єкт може зберігатися у вигляді точки, лінії або полігону.

На растровій карті властива об'єкту позиція визначається комірками з відповідними атрибутами. Кожна комірка характеризується розміром та значенням атрибутів, яке може бути відображене у вигляді значень кольорової моделі (спутниковий знімок), або значень атрибуту (GRID), на основі яких растр може бути класифікований за кількісними або якісними характеристиками.

Деякі системи також підтримують більш складні об'єкти, такі як регіони, кола, еліпси, дуги і криві.

### **Інформація про атрибути**

Атрибути описують специфічні характеристики об'єктів, які не є графічними елементами. Для прикладу, атрибутом, що характеризує дорогу може бути її назва або дата, коли її побудували. Зазвичай, атрибути зберігаються окремо від графічної частини карти за допомогою геореляційних структур даних. Атрибути є частиною опису векторних карт і рідко застосовуються разом із растровими зображеннями.

### **Інформація про відображення**

Інформацію, що відображається може описувати цифрова карта або накреслена (план). Загальна інформація на карті водображається за допомогою різних кольорів: озера, моря, вершини, гори, рівнин; ширини і типів ліній: назви доріг, вулиць, нанесення висотної поясності на планах тощо.

Інформація про відображення задає графічні стилі, за допомогою яких відображаються дані на карті. Основними графічними стилями є графічні перемінні, до яких відноситься колір, товщина, форма, орієнтація, символізація, підписи. Разом графічні перемінні утворюють способи умовних позначень на карті географічних об'єктів та формують легенду карти.

**GIS-формати** створені для стандартного кодування географічної інформації в файл. Файли цього формату створюються державними агентствами з картографування, комерційними організаціями або розробниками програмного забезпечення.

Ці формати є розповсюдженими, але існує програмне забезпечення, що не підтримує файли цього формату. В таких випадках шукають способи перетворення GIS-файлу або використовують інший файл. Майже кожен GIS має свій власний внутрішній формат файлу. Ці формати призначені для оптимального використання розробником і часто є запатентованими. Вони не призначені для використання поза системою, що їх розробляла.

### **Векторні формати**

Багато додатків GIS засновані на векторній технології через її широке застосування. Вони є складними, тому що інформацію можливо описати безліччю способів: для зберігання

координат, атрибутів, зв'язків, структур баз даних та відображення інформації. Найбільш поширеним форматом гепросторових даних є формат [Shapefile](https://en.wikipedia.org/wiki/Shapefile)<sup>40</sup> (shp). Цей формат розроблений компанією ESRI та є пропрієтарним форматом даних. Проте сьогодні він став дефакто стандартом обміну та публікації даних. Більшість програмних продуктів працюють з цим форматом, в тому числі ПЗ QGIS, яке є найбільш популярним OpenSource додатком для ГІС. Специфікація формату визначає набір файлів із загальним префіксом імені файлу, які зберігаються в тому ж каталозі. Три обов'язкових файли мають розширення файлів shp, .SHX і .dbf.

.shp - головний файл; містить набори геометричних об'єктів  
.shx - індексний файл; який використовується для зв'язку між файлами .dbf і .shp  
.dbf - атрибутивний файл; містить атрибути об'єктів, описаних в .shp файлі у форматі dBase IV.

На сучасному етапі більшість геоданих зберігаються в СУБД, які підтримують збереження просторових даних. Зараз до них відносяться Oracle, SQLserver, Postgres та ін. ПЗ ArcGis підтримує два типи баз даних: персональна та геобаз даних з відповідними форматами: mdb, gdt.

Значного поширення набули формати ГІС: [MapInfo TAB format](https://en.wikipedia.org/wiki/MapInfo_TAB_format)<sup>41</sup>, [Spatialite](https://en.wikipedia.org/wiki/Spatialite)<sup>42</sup>, [Geography Markup Language \(GML\)](https://en.wikipedia.org/wiki/Geography_Markup_Language)<sup>43</sup>. Також розповсюдженими в Україні є формати даних ПЗ Панорама – SXF, Digital – підтримка формату in4, який довгий час був обмінним файлом кадастрової інформації.

Файли AutoCAD не є ГІС форматом, проте сучасні професійні ГІС підтримують пряму конвертацію з формату DXF та навпаки. Доволі рідкими форматами в Україні є формати [National Transfer Format](https://en.wikipedia.org/wiki/National_Transfer_Format)<sup>44</sup> (NTF), [SOSI](https://en.wikipedia.org/wiki/SOSI)<sup>45</sup>, [Spatial Data File](https://en.wikipedia.org/wiki/Spatial_Data_File)<sup>46</sup>, [TIGER](https://en.wikipedia.org/wiki/TIGER)<sup>47</sup>, [Vector Product Format](https://en.wikipedia.org/wiki/Vector_Product_Format)<sup>48</sup> (VPF).

Для обміну просторовими даними також використовують формат ASCII, який у більшості випадків використовують для збереження растрових даних.

Для прикладу наведено найпоширеніші формати: [AutoCAD DXF](https://en.wikipedia.org/wiki/AutoCAD_DXF)<sup>49</sup>, [Cartesian coordinate system](https://en.wikipedia.org/wiki/Cartesian_coordinate_system)<sup>50</sup> (XYZ), [Digital Line Graph](https://en.wikipedia.org/wiki/Digital_line_graph)<sup>51</sup> (DLG), [Esri TIN](https://en.wikipedia.org/wiki/Esrri_TIN)<sup>52</sup>, [Geography Markup Language](https://en.wikipedia.org/wiki/Geography_Markup_Language)<sup>53</sup> (GML),

---

<sup>40</sup><https://en.wikipedia.org/wiki/Shapefile>

<sup>41</sup>[https://en.wikipedia.org/wiki/MapInfo\\_TAB\\_format](https://en.wikipedia.org/wiki/MapInfo_TAB_format)

<sup>42</sup><https://en.wikipedia.org/wiki/Spatialite>

<sup>43</sup>[https://en.wikipedia.org/wiki/Geography\\_Markup\\_Language](https://en.wikipedia.org/wiki/Geography_Markup_Language)

<sup>44</sup>[https://en.wikipedia.org/wiki/National\\_Transfer\\_Format](https://en.wikipedia.org/wiki/National_Transfer_Format)

<sup>45</sup><https://en.wikipedia.org/wiki/SOSI>

<sup>46</sup>[https://en.wikipedia.org/wiki/Spatial\\_Data\\_File](https://en.wikipedia.org/wiki/Spatial_Data_File)

<sup>47</sup><https://en.wikipedia.org/wiki/TIGER>

<sup>48</sup>[https://en.wikipedia.org/wiki/Vector\\_Product\\_Format](https://en.wikipedia.org/wiki/Vector_Product_Format)

<sup>49</sup>[https://en.wikipedia.org/wiki/AutoCAD\\_DXF](https://en.wikipedia.org/wiki/AutoCAD_DXF)

<sup>50</sup>[https://en.wikipedia.org/wiki/Cartesian\\_coordinate\\_system](https://en.wikipedia.org/wiki/Cartesian_coordinate_system)

<sup>51</sup>[https://en.wikipedia.org/wiki/Digital\\_line\\_graph](https://en.wikipedia.org/wiki/Digital_line_graph)

<sup>52</sup>[https://en.wikipedia.org/wiki/Esrri\\_TIN](https://en.wikipedia.org/wiki/Esrri_TIN)

<sup>53</sup>[https://en.wikipedia.org/wiki/Geography\\_Markup\\_Language](https://en.wikipedia.org/wiki/Geography_Markup_Language)

GeoJSON<sup>54</sup>, GeoMedia<sup>55</sup>, ISFC, Keyhole Markup Language<sup>56</sup> (KML), MapInfo TAB format<sup>57</sup>, National Transfer Format<sup>58</sup> (NTF), Spatialite<sup>59</sup>, Shapefile<sup>60</sup>, Simple Features<sup>61</sup>, SOSI<sup>62</sup>, Spatial Data File<sup>63</sup>, TIGER<sup>64</sup>, Vector Product Format<sup>65</sup> (VPF).

## Растрові формати

Растрові файли зазвичай використовуються для зберігання інформації зображення у вигляді сканованих паперових карт чи аерофотозйомки. Вони також використовуються для передачі даних дистанційного зондування Землі. На відміну від інших растрових файлів, які характеризуються розрізненням у розмірах точок на дюйм, розрізнення на зображеннях дистанційного зондування виражається в метрах на піксел растру, що вказує розмір земельної ділянки певної області. Прикладами растрових форматів є: ADRG, RPF, CADRG, CIB<sup>66</sup>, Digital raster graphic<sup>67</sup> (DRG), ECRG, ECW<sup>68</sup>, Esri grid<sup>69</sup>, GeoTIFF<sup>70</sup>, IMG – ERDAS IMAGINE<sup>71</sup>, JPEG2000<sup>72</sup>, MrSID<sup>73</sup>, netCDF<sup>74</sup>-CF.

Детальніше про формати та GIS-формат: [https://en.wikipedia.org/wiki/GIS\\_file\\_formats](https://en.wikipedia.org/wiki/GIS_file_formats)<sup>75</sup>

## KML - геоінформаційний формат Google Maps

KML (від англ. *Keyhole Markup Language*) - мова розмітки на основі XML для представлення тривимірних геопросторових даних в програмі «Google Earth».

Підмножина мови KML 2.0 може використовуватися і для відображення двовимірних карт в сервісі «Карти Google».

KML-файли зазвичай поширюються в ZIP-архіві: KMZ

Зміст файлів KML визначає один або кілька об'єктів для відображення в Google Earth. Цими об'єктами можуть бути:

---

<sup>54</sup><https://en.wikipedia.org/wiki/GeoJSON>

<sup>55</sup><https://en.wikipedia.org/wiki/GeoMedia>

<sup>56</sup>[https://en.wikipedia.org/wiki/Keyhole\\_Markup\\_Language](https://en.wikipedia.org/wiki/Keyhole_Markup_Language)

<sup>57</sup>[https://en.wikipedia.org/wiki/MapInfo\\_TAB\\_format](https://en.wikipedia.org/wiki/MapInfo_TAB_format)

<sup>58</sup>[https://en.wikipedia.org/wiki/National\\_Transfer\\_Format](https://en.wikipedia.org/wiki/National_Transfer_Format)

<sup>59</sup><https://en.wikipedia.org/wiki/Spatialite>

<sup>60</sup><https://en.wikipedia.org/wiki/Shapefile>

<sup>61</sup>[https://en.wikipedia.org/wiki/Simple\\_Features](https://en.wikipedia.org/wiki/Simple_Features)

<sup>62</sup><https://en.wikipedia.org/wiki/SOSI>

<sup>63</sup>[https://en.wikipedia.org/wiki/Spatial\\_Data\\_File](https://en.wikipedia.org/wiki/Spatial_Data_File)

<sup>64</sup><https://en.wikipedia.org/wiki/TIGER>

<sup>65</sup>[https://en.wikipedia.org/wiki/Vector\\_Product\\_Format](https://en.wikipedia.org/wiki/Vector_Product_Format)

<sup>66</sup>[https://en.wikipedia.org/wiki/Controlled\\_Image\\_Base](https://en.wikipedia.org/wiki/Controlled_Image_Base)

<sup>67</sup>[https://en.wikipedia.org/wiki/Digital\\_raster\\_graphic](https://en.wikipedia.org/wiki/Digital_raster_graphic)

<sup>68</sup>[https://en.wikipedia.org/wiki/ECW\\_\(file\\_format\)](https://en.wikipedia.org/wiki/ECW_(file_format))

<sup>69</sup>[https://en.wikipedia.org/wiki/Esri\\_grid](https://en.wikipedia.org/wiki/Esri_grid)

<sup>70</sup><https://en.wikipedia.org/wiki/GeoTIFF>

<sup>71</sup>[https://en.wikipedia.org/wiki/ERDAS\\_IMAGINE](https://en.wikipedia.org/wiki/ERDAS_IMAGINE)

<sup>72</sup><https://en.wikipedia.org/wiki/JPEG2000>

<sup>73</sup><https://en.wikipedia.org/wiki/MrSID>

<sup>74</sup><https://en.wikipedia.org/wiki/NetCDF>

<sup>75</sup>[https://en.wikipedia.org/wiki/GIS\\_file\\_formats](https://en.wikipedia.org/wiki/GIS_file_formats)

- **Позначки на карті.** Задається знаковий символ, який відображається в певному місці на карті, і його параметри: величина, текст і колір, величина написів, яка відображається біля знаку. Є можливість створити позначку без знаку. Наприклад, поставити номери будинків району. Для цього потрібно просто в стилі створити порожнє посилання на значок:

<Icon> <href> </href> </Icon>

- **Багатокутник або набір ліній.** Визначається колір ліній і колір підпису.
- **Зображення.** Визначається положення зображення на поверхні Землі, а також його масштаб. Також можна розмістити зображення на екрані, що не переміщається разом з картою - наприклад, логотип.
- **Тривимірна модель.** Версія мови KML 2.1 (яка відповідає четвертій версії програми Google Earth) дозволяє підключати опис тривимірних об'єктів (наприклад, будівель і споруд). Тривимірну модель можна задати двома способами: завданням висоти плоских фігур (витягуванням) і посиланням на повноцінну модель в форматі Collada.

Крім згаданих характеристик, для кожного об'єкта також задаються основні геоінформаційні властивості (географічна широта і довгота, а також висота над рівнем моря або над рівнем поверхні Землі). Може бути наведено короткий опис об'єкта (який в подальшому відображається в підказці за запитом користувача). Може бути зазначений рекомендований ракурс спостереження зазначеного на карті місця, тобто рекомендована висота, азимут і кут нахилу «віртуальної камери».

Об'єкти всередині KML-файлу можна організувати в ієрархічні структури папок і підпапок, щоб було зручніше спільно ввімкнути або вимкнути відображення логічно взаємопов'язаних груп об'єктів.

KML-файл може містити (в URL-формі) посилання на інші файли в форматі KML або KMZ, розташовані де-небудь в мережі, і задавати умови і регулярність завантаження і відображення даних з цих зовнішніх джерел. Таке мережеве посилання також видно в якості своєрідної підпапки.

Мова KML багато в чому наслідує структуру мови GML - географічного мови розмітки.

Приклад KML-розмітки:

```

1  <?xml version="1.0" encoding="UTF-8"?>
2  <kml xmlns="http://earth.google.com/kml/2.1">
3  <Placemark>
4      <name>Київ</name>
5      <description><![CDATA[<p>Київ, Київська область, Україна.</p>
6          Столиця України. Місто побудоване на березі річки Дніпро.]]>
7      </description>
8      <LookAt id="khLookAt540_copy0">
9          <longitude>38.0576198113139</longitude>
10         <latitude>44.56963150481845</latitude>
11         <altitude>0</altitude>
12         <range>14693.40972993507</range>
13         <tilt>49.10268313434742</tilt>
14         <heading>37.85562764777833</heading>
15     </LookAt>
16     <Style>
17         <IconStyle>
18             <scale>0.9</scale>
19             <Icon>
20                 <href>root://icons/palette-4.png</href>
21                 <x>32</x>
22                 <y>128</y>
23                 <w>32</w>
24                 <h>32</h>
25             </Icon>
26         </IconStyle>
27         <LabelStyle>
28             <scale>0.9</scale>
29         </LabelStyle>
30     </Style>
31     <Point id="khPoint541_copy0">
32         <coordinates>
33             38.06284424434902,44.56842733252498,0
34         </coordinates>
35     </Point>
36 </Placemark>
37 </kml>

```

## GeoJSON та TopoJSON

GeoJSON - відкритий формат, призначений для зберігання графічних структур даних, заснований на форматі JSON.



Формат може зберігати примітивні типи для опису географічних об'єктів, такі як: точки (адреси та місця розташування), лінії (вулиці, шосе, кордони), полігони (країни, штати, ділянки землі). Також можуть зберігатися так звані мультитипи, які представляють собою об'єднання декількох примітивних типів.

Формат GeoJSON відрізняється від інших стандартів ГІС тим, що він був написаний і підтримується не певною організацією зі стандартизації, а за допомогою робочої групи розробників.

Подальшим розвитком GeoJSON є TopoJSON, розширення GeoJSON, яке кодує геопросторову топологію, і, як правило, забезпечує менший розмір файлів. Сьогодні ці формати підтримуються більшістю ГІС, а також безпосередньо застосовуються при публікації інтерактивних картографічних даних в веб, зокрема засобами бібліотек OpenLayers та Leaflet.

```
1 { "type": "FeatureCollection",
2   "features": [
3     { "type": "Feature",
4       "geometry": { "type": "Point", "coordinates": [102.0, 0.5] },
5       "properties": { "prop0": "value0" }
6     },
7     { "type": "Feature",
8       "geometry": {
9         "type": "LineString",
10        "coordinates": [
11          [102.0, 0.0], [103.0, 1.0], [104.0, 0.0], [105.0, 1.0]
12        ]
13      },
14      "properties": {
15        "prop0": "value0",
16        "prop1": 0.0
17      }
18    },
19    { "type": "Feature",
20      "geometry": {
21        "type": "Polygon",
22        "coordinates": [
23          [ [100.0, 0.0], [101.0, 0.0], [101.0, 1.0],
24            [100.0, 1.0], [100.0, 0.0] ]
25        ]
26      },
27      "properties": {
28        "prop0": "value0",
29        "prop1": { "this": "that" }
30      }
31    }
32  ]
33 }
```

```
31         }  
32     ]  
33 }
```