

---

# Лекція 3

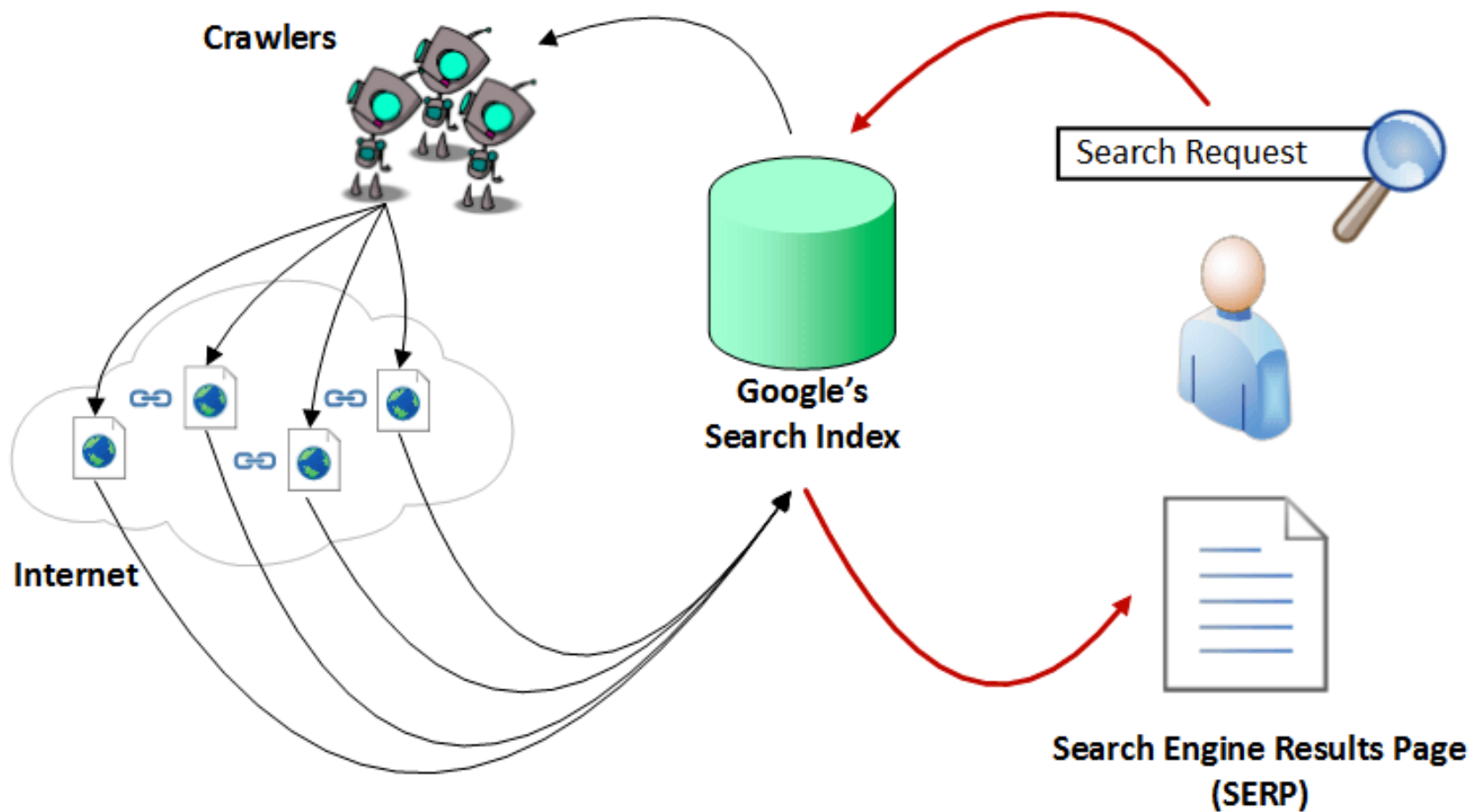
## Моделі пошуку. Метакраулери.

Курс «Робота з інформаційними технологіями»

Яворський В.А.

Листопад 2017 рік

# Схема процесу пошуку



# Класичні моделі пошуку

---

- *Булева модель*
- *Векторно-просторова модель*
- *Ймовірнісна модель*

Припущення:

*Розгляд документів як множини окремих слів,  
незалежних одне від одного - концепція «Bag of Words».*

# Булева модель пошуку

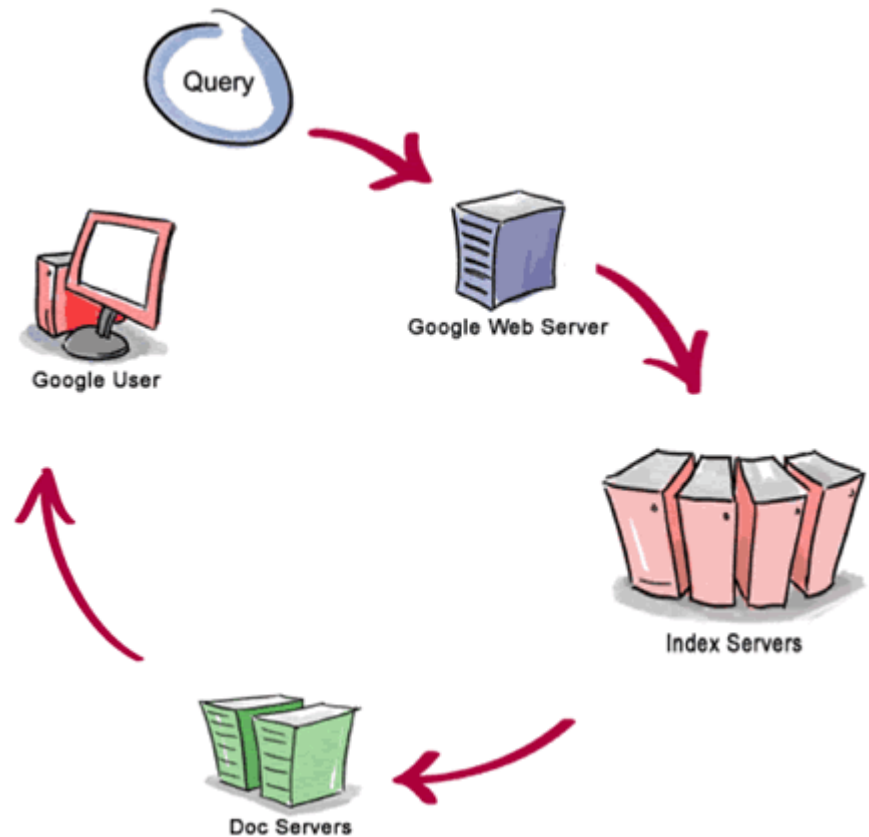
**Запит** - логічне вираження з операторами (AND, OR, NOT)

**Склад таблиць з інвертованими списками:**

- Тексти;
- Індеси на тексти;
- Словник унікальних слів;
- інверсна, що містить списки номерів документів, які відповідають певним словами.

**Процес пошуку інформації в ІПС з інвертованими списками:**

- Звернення до словника унікальних слів;
- Звернення до інверсної таблиці;
- Звернення до покажчиків на тексти;
- Звернення до текстової таблиці.



# Векторно-просторова модель пошуку

**Документ** - описується вектором в деякому евклідовому просторі термінів. Кожному терміну зіставляється вага, яка характеризується частотою, місцем розташування, тематикою та інше.

## **Вага:**

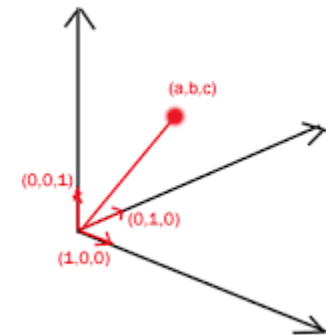
множимо на частоту появи терміну в документі  
множимо на обернену величину кількості документів масиву, де зустрічається термін

**Запит** - також вектор в евклідовому просторі.

Близькість запиту документу – скалярний добуток.

Модель забезпечує:

- Обробку запитів без логічних обмежень їх довж
- Простоту реалізації режиму ПОСК подібних документів;
- Збереження результатів пошуку з можливістю виконання уточнюючого пошуку.



# Ймовірнісна модель пошуку

---

- Беремо апріорні оцінки ймовірності того, що документ є релевантним - виходячи зі складу термінів.
- Отримуємо експертні оцінки користувачів, які визнають документ релевантним або нерелевантним.
- На кожному етапі ітерації, завдяки режиму зворотного зв'язку, визначаються документи, зазначених користувачем як такі що задовольняють його інформаційні потреби (є пертинентними).

*Модель добре визначає новий спам по множині документів, які визначені користувачами як спам.*

# Недоліки класичних моделей

---

- **Булева модель** - невисока ефективність пошуку, жорсткий набір операторів, неможливість рангування.
- **Векторно-просторова модель** - пов'язана з розрахунком масивів високої розмірності, малоприматна для обробки великих масивів даних.
- **Ймовірнісна модель** - має низьку обчислювальну масштабованість, пов'язана з необхідністю постійного навчання системи.

# Google використовує понад 200 факторів для визначення релевантності





# Google- історія



[« Особиста інформація та конфіденційність](#)

## Інформаційна панель Google

Щомісяця надсилати мені нагадування про перевірку активності в обліковому записі. [i](#)

[Розгорнути все](#)

### Обліковий запис

Назва	Електронна адреса облікового запису ...	Керувати обліковим записом
Владимир Яворский	vlayar@gmail.com	<a href="#">Змінити пароль</a> <a href="#">Пов'язані додатки та сайти</a>

### Android

Пристрої	Керувати активними пристроями
<b>2</b>	<a href="#">Пристрої в магазині Google Play</a>

### Blogger

Назва	Редагувати профіль Blogger
Владимир Яворский	<a href="#">Керувати блогами</a>

### Gmail

Бесіди	Остання бесіда	Налаштування
<b>1 458</b>	<a href="#">Последняя цена на товары со склада ...</a>	<a href="#">Конфіденційність і безпека</a>

### Google+

Оцінки +1	Редагувати профіль
<b>2</b>	<a href="#">Редагувати кола</a> <a href="#">Публікації</a>

[https://www.google.com/  
settings/activity/](https://www.google.com/settings/activity/)

[https://myaccount.google.  
com/dashboard](https://myaccount.google.com/dashboard)

# «Чорні» методи пошукового спаму

---

- Сторінка з великою частотою пошукового слова або словосполучення
- Прихований, дрібний і невидимий текст  
приховані малюнки
- Автоматична переадресація
- Неадекватні ключові слова і опис
- Велика кількість однакових сторінок однієї тематики

# «Чорні» методи пошукового спаму

---

- **flood** («затоплення» пошукової системи) - індексування сторінки під різними мережевими іменами
- Перевищення числа сторінок в заявці на індексацію пошуковими краулерами
- **Дорвей** (doorway pages) - сторінки, що містять розрізнені набори ключових слів на найрізноманітніші теми, створені для пошукових роботів.

# «Чорні» методи пошукового спаму

---

- **Своппінг** (swapping) - оптимізація сторінок для досягнення верхніх позицій результатах пошуку з наступною заміною змісту
- **Клоакинг** (cloaking) - програмне забезпечення на сервері здатне розпізнавати роботів пошукових систем і підставляти їм не той зміст сторінок, яке побачать відвідувачі

*не помітить пошуковий робот - помітять відвідувачі! І поскаржаться адмінам ...*

# «Сірі» методи пошукового спаму

---

- Сторінки з швидким оновленням
- Дзеркала сайту
- «Пошуково-активні» вхідні сторінки
- Мережі обміну посиланнями
- Багаторівневий маркетинг (MLM-технології)
- «Стратегічний» спам для індексу цитування

# Оновлення алгоритмів пошуку Google 2011-2015 роки

---

## **Panda (лютий 2011)**

*значне поліпшення алгоритму пошуку, яке спрямоване на підвищення якості контенту веб-сайтів. Оригінальні сайти з авторським контентом в пошуковій системі повинні зайняти місце вище, ніж сторінки з низькою якістю, що повторюють те, що вже і так відомо або ж є копіями інших сайтів.*

- вміст на сторінці повинен мати істотний обсяг >1500 слів;
- інформація, представлена на сайті повинна бути оригінальною. Якщо ви просто копіюєте вміст інших веб-ресурсів - Google покарає;
- оригінальність - для успішного просування контент має бути те, чого немає на інших сайтах;
- текст сайту повинен бути орфографічно і граматично правильним, зміст повинен відповідати описаним стандартам.

# Оновлення алгоритмів пошуку Google

---

## **Page Layout (січень 2012)**

*Покарання сайтів, які використовують занадто багато реклами у верхній частині сторінки або роблять її надмірно агресивною, що відволікає від основного змісту (користувачам складно знайти потрібну інформацію і доводилося довго прокручувати сторінку вниз, велика кількість реклами заважає зручності засвоєння інформації).*

## **Penguin (березень 2012)**

*Боротьба з пошуковим спамом. Сайти, які використовували спам-методи, були значно знижені в рейтингу або зовсім вилучені з нього. Здатність аналізувати кількість посилань.*

## **Pirate (серпень 2012)**

*Зниження рейтингу сайтів, які порушують авторські права та інтелектуальну власність. Для оцінки цих порушень, Google використовує систему запитів про порушення авторських прав, засновану на Digital Millenium Copyright Act.*

# Оновлення алгоритмів пошуку Google

---

## **Exact Match Domain (EMD, вересень 2012)**

*Боротьба з доменами, схожими на MFA.*

*MFA (made-for-adsense) - домен, який створений спеціально для контекстно-медійної системи Google. Зазвичай такий домен призначений для якогось одного запиту (або сімейства запитів) і на ньому встановлений Google Adsense. Користувач, який потрапив на цей домен, не бачить нічого, крім реклами і або закриває сайт, або переходить далі по контекстному оголошенню.*

## **Payday Loan (червень 2013)**

*Зменшення кількості сторінок, які містять переспамлені запити.*

*Приклад: вам потрібно купити двері. На запит Google видасть фотографії дверей. З них: 2-3 сторінки, де безпосередньо можна купити двері, 3-4 сайту компаній-виробників дверей і 2-3 сайту про те, як вибрати і поміняти двері. Якби не було оновлення Payday Loan, ви б побачили 15-20 запитів на одну тематику (наприклад, де купити двері).*



# Оновлення алгоритмів пошуку Google

---

## **Hummingbird (вересень 2013)**

*Аби повертати точні відповіді на запити із ключовими словами, Google інтерпретує наміри і контекст пошуку. Мета полягає в тому, щоб зрозуміти сенс пошукового запиту користувача і повертати відповідні результати. Це означає, що точні співпадіння ключових слів стають менш важливими на користь пошуку наміри.*

*Приклад: якщо ви вводите запит «погода», то навряд чи очікуєте отримати повне пояснення самого терміна, натомість отримаєте опис погодних умов.*

## **Pigeon (липень 2014)**

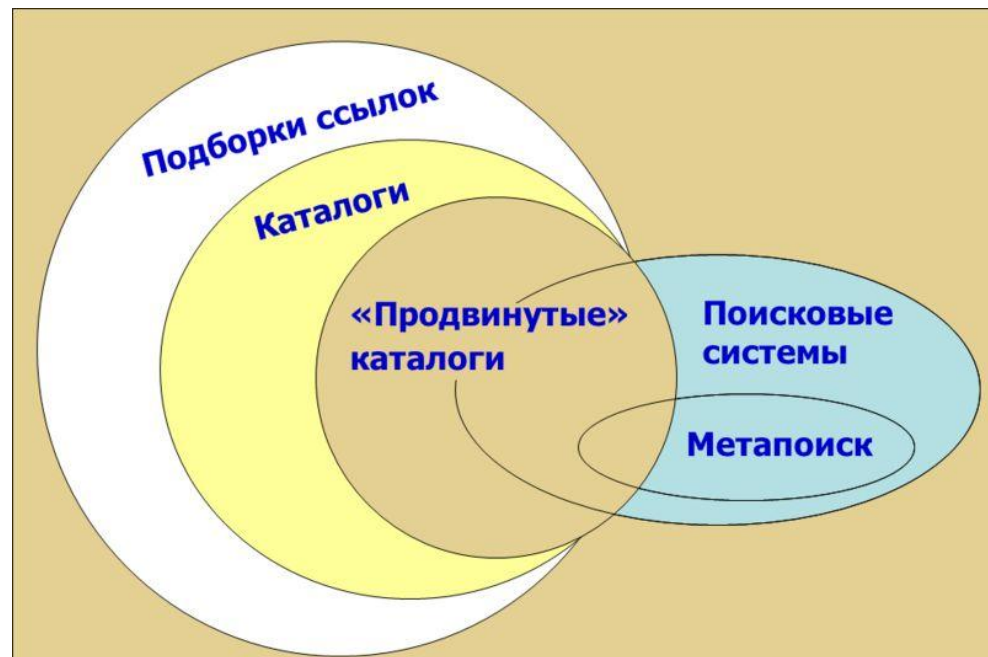
*Геозалежний пошук. Відстань і місце розташування користувача є ключовими параметрами ранжування, аби забезпечити точність результату.*

## **Mobilegeddon (квітень 2015)**

*Мобільний пошук. Google дає перевагу сторінкам, дружнім до мобільних пристроїв.*

# Типи пошукових систем

- Системи з пошуковими роботами
- Каталоги ресурсів (керовані людиною)
- Довідникові ресурси
- Локальні програми для пошуку в інтернеті
- Гібридні системи
- Метасистеми



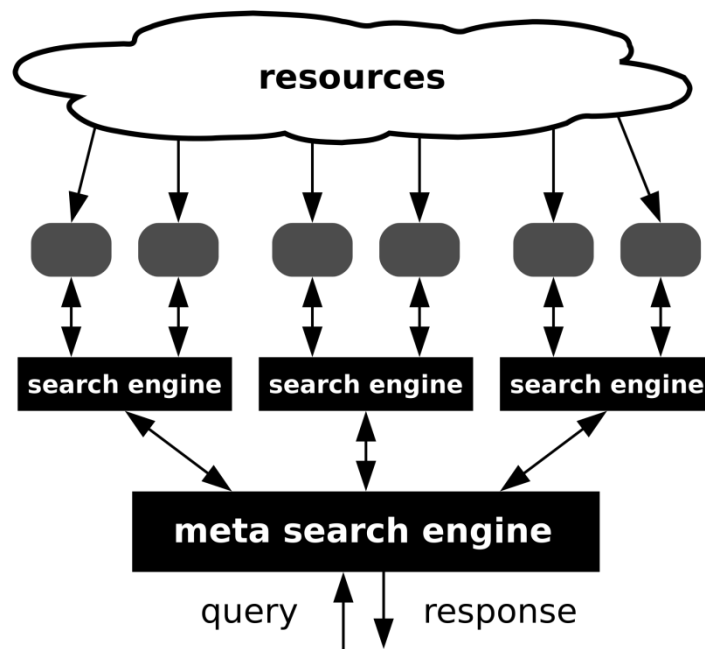
# Веб-каталоги

---

- Принципово: в наповненні ресурсами беруть участь люди, а не автоматичні пошукові програми.
- Рекомендовано для першого знайомства з предметною областю
- Рекомендовано для пошуку по нечітким запитам
- Недоліки:
  - Слабка оперативність
  - Розмір баз
  - Відсутня єдина класифікація ресурсів і чіткі критерії віднесення до категорій

# Метапошукові системи (searchbots, metacrawlers)

Замість свого пошуку та індексування інформації, метасистеми надсилають запити багатьом пошуковим ресурсам.



# Метапошукові системи

---

<https://www.startpage.com/>

<https://www.ixquick.com/>

- Цільові ресурси з кількох пошукових систем
- анонімність серфінгу - комп'ютери клієнтів не зв'язуються з сайтами безпосередньо.
- Політика конфіденційності Startpage - IP-адреси, відомості про відвідувані портали, ключові терміни, куки не зберігаються. За це нагороджений першим Європейським Знаком Конфіденційності (EuroPriSe) 14 липня 2008 року

start  
page by ixquick

# Етапи метапошуку

---

- ✓ Не тільки пара сторінок з результатами різних пошукових систем
- ✓ Технологія кластеризації результатів пошуку

## Етапи:

1. Пошук веб-сторінок по запиті;
2. Аналіз знайдених сторінок, знаходить додаткові ключові слова, які зустрічаються разом з термінами у запиті;
3. Підмножини сторінок в результаті аналізу оголошуються кластерами;
4. Визначення релевантності посилань і їх позиції в результатах кожного окремого кластера, посилання в межах кластера мають вищу ціну;
5. Рангування виділяє в кластерах корисні ресурси, яким при звичайному пошуку мало що «світить»

# Nigma.ru

---

- Власний алгоритм кластеризації результатів
- Врахування специфіки російськомовних запитів
  - Пошук по різним словоформам, синонімам та узагальненим поняттям
  - Потужна система виправлення орфографічних помилок
- Автоматичне доповнення введених запитів та переклад
- Рішення математичних, хімічних рівнянь
- Структурована інформація по запиту



**NiGMA.RU**

интеллектуальная поисковая система



со **всеми** словами:  с **точной** фразой:

с **любым** из слов:  **без слов:**

на **сайте:**  *Пример:*

[в Киеве \(Укра...](#)  [Поисковики](#)  [Язык](#)

- Yandex
- Google
- Rambler
- Bing
- Yahoo
- Altavista
- Nigma

адаптация

В найденном  в Москве

50 результатов.

- [Официальный сайт - Адапт...](#)  
«Адаптация» — казахстанская русскоязы...  
[Читать дальше на Википедии →](#)  
Сайт (рус.): [ermen.antimusic.ru](#)  
Понравился поиск? Сделай Нигма.рф [поиски](#)
- [Адаптация все серии \(23.02.201...](#)  
Хочу сегодня поделиться своим мнением  
[Найти слова](#) | [livefilm.info/ser/67927-adapta](#)
- [Адаптация 11, 12 серия \(2017\) с...](#)  
Смотреть онлайн **Адаптация** сериал ТНТ.  
[Найти слова](#) | [kinoclips.tv/news/adaptacija](#)
- [Адаптация 12 серия \(сериал 201...](#)  
**Адаптация** сериал ТНТ смотреть онлайн.  
[Найти слова](#) | [kinoclips.tv/news/adaptacija](#)
- [Русский сериал Адаптация 1, 2...](#)  
Американским властям следует оперативн...  
разведки свидетельствуют...  
[Найти слова](#) | [minizal.su/serialw3815-adapt](#)
- [Джон Уиндем | Адаптация](#)

адаптация

- [адаптация](#) Физиологическая адаптация [ru.wikipedia.org/...](#)
- [адаптация](#) детей к детскому саду [7va.ru/...](#)
- [адаптация](#) персонала [ru.wikipedia.org/...](#)
- [адаптация](#) это [ru.wikipedia.org/...](#)
- [адаптация](#) ребенка в детском саду [volgo-mame.ru/...](#)

Нигма.Справка

**Биологическая адаптация** (от лат. adaptatio — приспособление) — приспособление организма ко внешним условиям в процессе эволюции, включая морфофизиологическую и поведенческую составляющие. Адаптация может обеспечивать выживаемость в условиях конкретного местообитания, устойчивость к воздействию факторов абиотического и биологического характера, а также успех в конкуренции с

Фильтр

Как это помогает искать?

- адаптация это что такое
- адаптация определение
- сериал адаптация
- смотреть онлайн
- тнт все серии
- процесс
- понятие
- сериал адаптация 2017
- тнт смотреть онлайн
- адаптация все серии
- смотреть онлайн
- в хорошем качестве
- адаптации
- Русскоязычные сайты

Фильтровать

Со всеми:  сбросить  выбрать  исключить

Вы не авторизованы

Ваше имя

[Регистрация](#)

[Напомнить?](#)

[Новости](#) | [Форум](#)

название	Адаптация
автор	Джон Уиндем

[Больше книг](#)

[Найти слова](#) | [lib.ru/NOFANT/UINDEM/26-23](#)




C<sub>2</sub>H<sub>5</sub>OH + CH<sub>3</sub>COOH

В найденном в Киеве (Украи... Поисковики Язык Сортировка Настройки

54 млн. результатов.

**Реакция:**  
(видео: [почему происходят химические реакции](#))

**Реакция образования сложных эфиров**



**Условия:**  
В присутствии сильной кислоты. Реакция идёт до конца при наличии водоотнимающего агента (напр.

[Больше реакций](#)

**Реакция:**  
CH3COOH + C2H5OH = CH3COOC2H5 + H2O

**Реакция:**  
CH3COOC2H5 + H2O = CH3COOH + C2H5OH

**Условия:** В присутствии H<sup>+</sup>

[Список решаемых задач](#)

1. Ответы@Mail.Ru: [C<sub>2</sub>H<sub>5</sub>OH+CH<sub>3</sub>COOH=](#) помогите пожалуйста))

C<sub>2</sub>H<sub>5</sub>OH+CH<sub>3</sub>COOH = CH<sub>3</sub>COOC<sub>2</sub>H<sub>5</sub> (этилацетат) + H<sub>2</sub>O. Следует указать что условием реакции этерификации является кислая среда! ...

[Найти слова](#) | [qvet.mail.ru/question/75713978](#) 951 6

log(x+10)\*(20-x)=0

В найденном в Киеве (Украи... Поисковики Язык Сортировка Настройки

125 млн. результатов.

Дано:

$$\ln(x + 10) \cdot (20 - x) = 0$$

Скрыть Решение

1 ОДЗ уравнения:

$$x \in (-10, \infty)$$

2 Делаем преобразование левой части уравнения:

$$\ln(x + 10) \cdot (20 - x) = -(x - 20) \cdot \ln(x + 10)$$

3 Уравнение после преобразования:

$$-(x - 20) \cdot \ln(x + 10) = 0$$

4 Решаем уравнение:

$$x = 20$$

5 Решаем уравнение:

$$-x = 9$$

6 Возможные решения:

$$-9;$$

$$20;$$

Ответ: (Решение уравнения с учётом ОДЗ)


$$x = -9,$$

$$x = 20$$

[Что это такое?](#) [Список решаемых задач](#) [Пожаловаться](#)

# Табличный Nigma-поиск

[Интернет](#) [Картинки](#) [Книги](#) [Музыка](#) [Математика](#) [Мини-игры](#)

[Расширенный поиск](#) 

Категория : Радиостанции, Класс : Радиостанции

Найти!

В найденном

в Киеве (Украи...

Поисковики

Язык

Сортировка

Настройки

92 статьи.

Регион поиска : **Киев** → [искать во всех регионах](#)

Исходный запрос : **радиостанции**

Показана таблица : **Радиостанции** → [искать по всем категориям](#)



Для перехода к другим таблицам используйте фильтр слева.  
Например, [Телеканалы](#) из категории **радиостанции**.

[Распечатать](#) | [Скачать таблицу](#) | [Ещё колонки](#) ↓

Статья	On-line трансляция	×	Веб-сайт	×	Владелец	×	Время вещания	×
<a href="#">100.9 FM</a>					«Румедиа»			
<a href="#">Best-FM</a>			<a href="#">best-fm.ru</a>		News Media Radio Group			
<a href="#">Business FM</a>			<a href="#">www.radio.businessfm.ru</a>		Аркадий Гайдамак			
<a href="#">Deutsche Welle</a>	<a href="#">аудио</a> <a href="#">видео</a>		<a href="#">www.dw-world.de</a>		ARD			
<a href="#">DFM</a>			<a href="#">dfm.ru</a>		«Русская Медиагруппа»			
<a href="#">Heart FM</a>			<a href="#">heartfm.ru</a>		ФГУП ГТРК Алтай			
<a href="#">Love Radio</a>			<a href="#">http://www.loveradio.ru/</a>		«Медиа холд»			
<a href="#">MAXIMUM</a>	<a href="#">http://www.maximum.ru/online/</a>		<a href="#">http://maximum.ru/</a>		Холдинг «Русская Медиагруппа»			
<a href="#">Megapolis FM</a>			<a href="#">megapolisfm.ru</a>					

Результаты поиска ограничены. [Отменить ограничения](#).

# Табличний Нігма-пошук

[Интернет](#) [Картинки](#) [Книги](#) [Музыка](#) [Математика](#) [Мини-игры](#)

википедия торрент

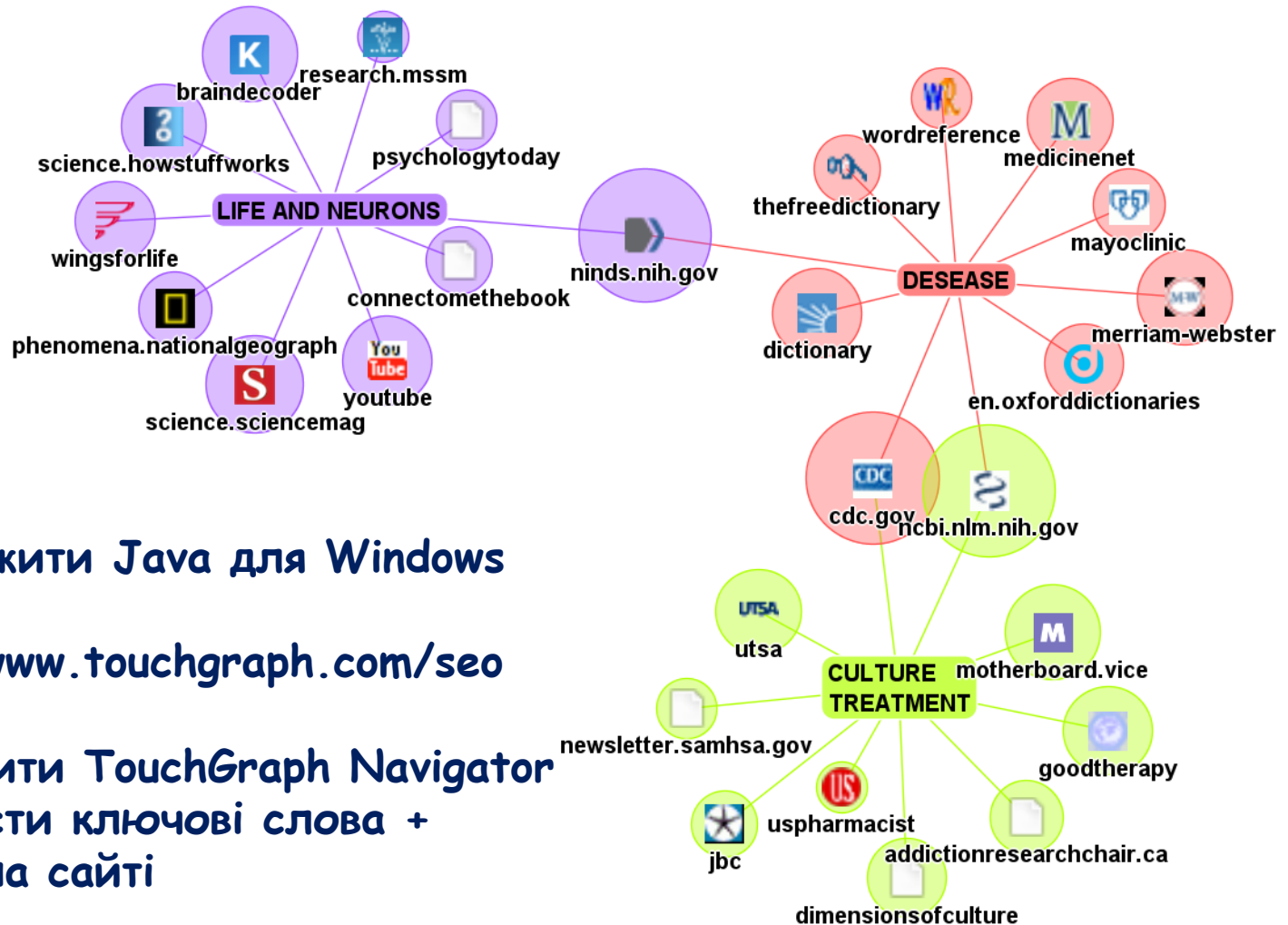
В найденном  в Киеве (Украи...

15 торрентов.

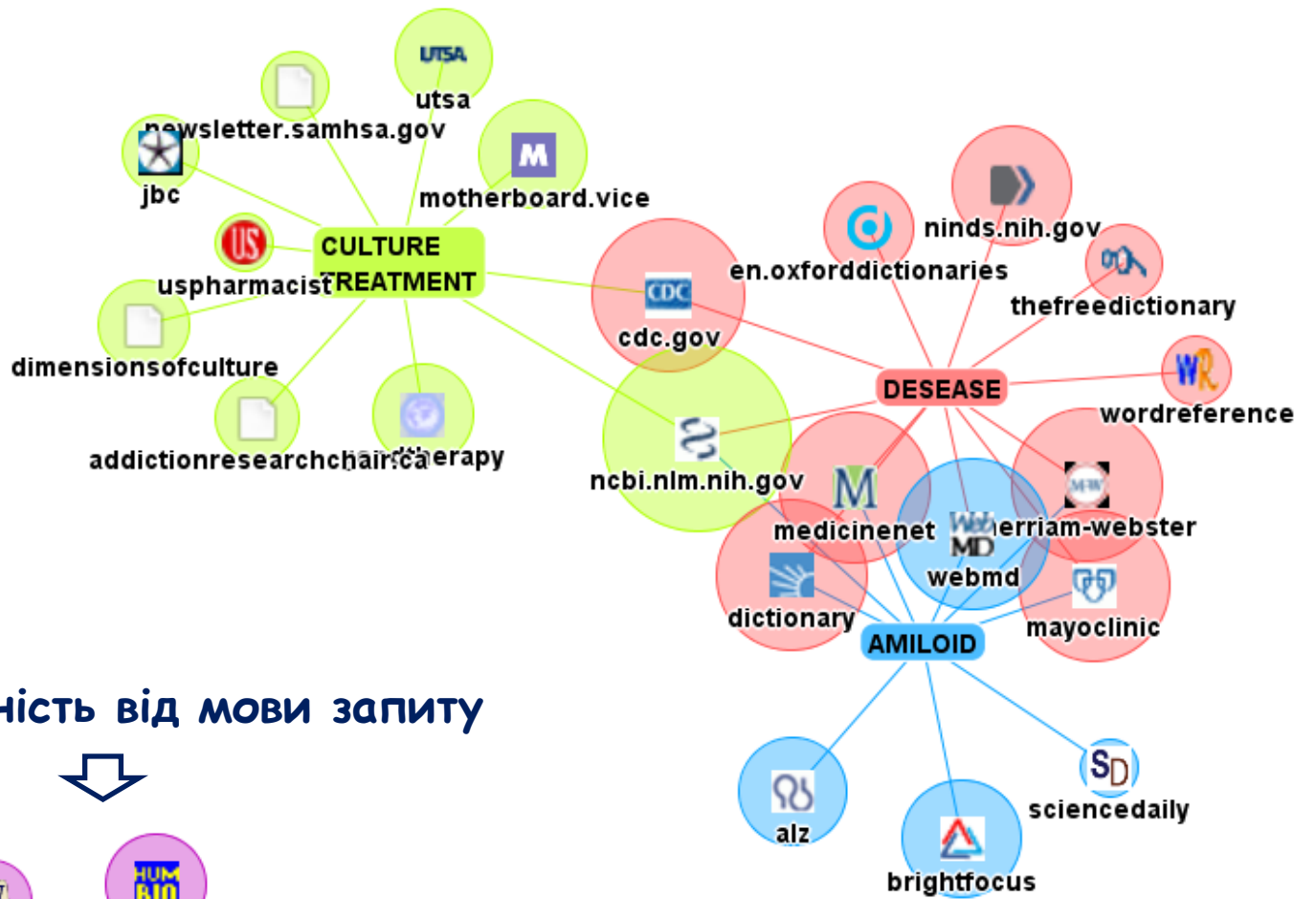
Показаны **только** торренты [показать все результаты](#)

Название	Размер	↑ Раздают	↓ Качают
<a href="#">Русская Википедия на 03.07.2010</a> Русская <b>Википедия</b> на 03.07.2010 (разбитая на 2 файла для запуска при недостатке оперативы) <a href="#">torrentino.com→</a>	799.1 Мб	0	0
<a href="#">Русская Википедия от 27.09.2010</a> Русская <b>Википедия</b> от 27.09.2010 [WM] <a href="#">torrentino.com→</a>	880.1 Мб	0	0
<a href="#">Русская Википедия Оффлайн - от 08.03.2012</a> Russian Wikipedia Offline / Русская <b>Википедия</b> Оффлайн - от 08.03.2012 [2012, RU] - WikiTaxi ver. 1.3.0 [2012] <a href="#">torrentino.com→</a>	2.9 Гб	0	0
<a href="#">Русская Википедия Оффлайн / Russian Wikipedia Offline от 2012.08.11 для DictViewer 2.0</a> Русская <b>Википедия</b> Оффлайн / Russian Wikipedia Offline от 2012.08.11 для DictViewer 2.0 [WM 6.x, RUS + ENG] <a href="#">torrentino.com→</a>	2.9 Гб	0	0
<a href="#">Русская Википедия Оффлайн / Russian Wikipedia Offline дамп ZD, от 2013.02.06 для Dictan, Dict</a> Русская <b>Википедия</b> Оффлайн / Russian Wikipedia Offline дамп ZD, от 2013.02.06 для Dictan, Dict [Android, WM 6.x, Windows 7, XP, RUS + ENG] <a href="#">torrentino.com→</a>	1.8 Гб	3	0
<a href="#">Украинская Википедия Оффлайн - от 07.12.2012 - WikiTaxi ver. 1.3.0</a> Ukrainian Wikipedia Offline / Украинская <b>Википедия</b> Оффлайн - от 07.12.2012 - WikiTaxi ver. 1.3.0 [2012, UK] <a href="#">torrentino.com→</a>	1.2 Гб	0	0
<a href="#">[Справочник] Русская Википедия от 01.06.11 [WM2003-6.x, RUS]</a>	1.2 Гб	0	0

# Візуальні пошукові системи

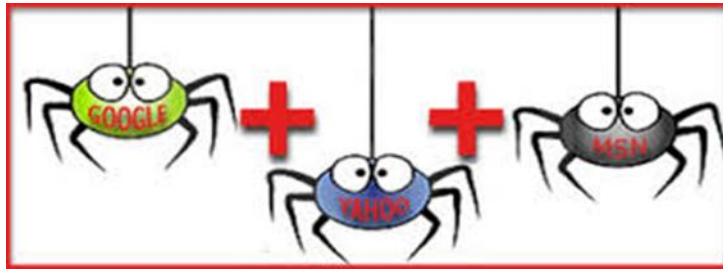


- Завантажити Java для Windows
- <http://www.touchgraph.com/se0>
- Встановити TouchGraph Navigator або ввести ключові слова + запуск на сайті



Висока залежність від мови запиту





## Завдання:

- 1. Створити інтелект-карту сайтів, які містять інформацію по Вашій темі дисертаційної роботи**
  - Підготувати ключові слова
  - Налаштувати візуальну пошукову метасистему
  - Провести пошук по кільком варіантам, обрати найкраще
  - Зберегти зображення та надіслати [jva@biph.kiev.ua](mailto:jva@biph.kiev.ua)