

Architecture of Information Systems

Search Engine

John Samuel
CPE Lyon

Year: 2017-2018

Email: john(dot)samuel(at)cpe(dot)fr



Search

Outline: Search engine

- Frontend development
- Backend development
- Application programming interface

Search Engine

Search

Search Engine

Search

Time

Location

...

Target Audience

- Regular users
- Domain Experts



Front-end development

Search interface

- One-box search
- Advanced search (filters)

Personalized user experience

Interface: Simple search (One box)

Queries

- Keywords
- Natural language queries

Search Results

- Links
- Structured response
- Natural language response

Léonard de Vinci — Wikipédia

https://fr.wikipedia.org/wiki/Léonard_de_Vinci ▼

Autoportrait de **Léonard de Vinci** réalisé entre 1512 et 1515, 33 × 21,6 cm , bibliothèque royale de Turin. Naissance. 15 avril 1452 · Vinci, Drapeau de la ...

[Homme de Vitruve](#) · [Liste des peintures de ...](#) · [La Cène \(Léonard de Vinci\)](#) · [Méduse](#)



Plus d'images

Léonard de Vinci

Peintre

Léonard de Vinci, né à Vinci le 15 avril 1452 et mort à Amboise le 2 mai 1519, est un peintre florentin et un homme d'esprit universel, à la fois artiste, organisateur de spectacles et de fêtes, ... [Wikipédia](#)

Date et lieu de naissance : 15 avril 1452, Anchiano, Italie

Date et lieu de décès : 2 mai 1519, Château du Clos Lucé, Amboise

Exposée : Musée du Louvre, Musée des Offices, PLUS ▼

Périodes : Haute Renaissance, Première Renaissance, Renaissance, Renaissance italienne, École florentine

Lieu d'inhumation : Chapelle Saint-Hubert, Amboise

Maître : Andrea del Verrocchio

Élève : Salai, Francesco Melzi...

Advanced search (filters)

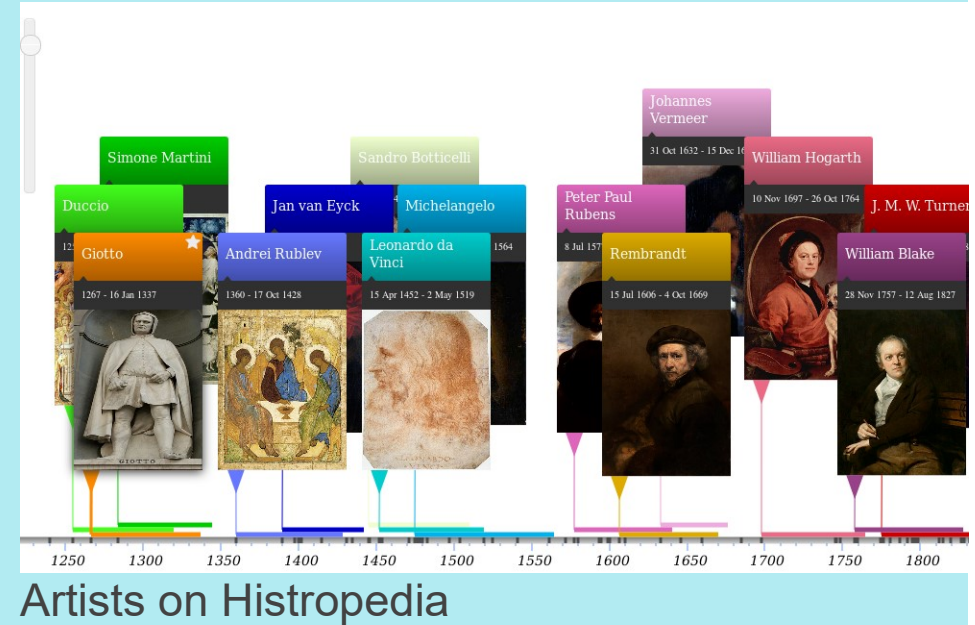
Filter search results (Multiple boxes)

Search

Time

Location

...



Advanced search (filters)



Location of Archaeological sites (Wikidata)

Advanced search (filters)

Why filters?

- Reduce information overload
- Precise queries
- Interactive search

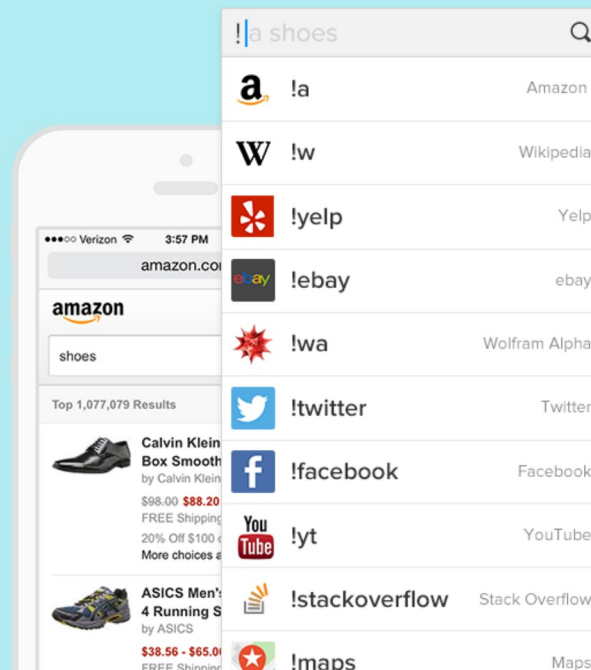
Advanced Search in one-box

Operators

- AND
- OR
- NOT

Advanced Search in one-box

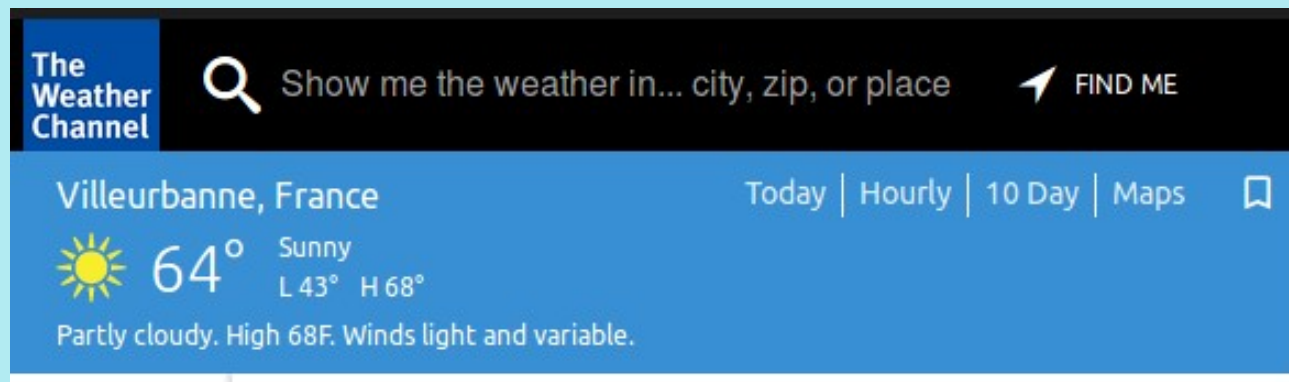
Mnemonics



Bangs (DuckDuckGo)

Personalized user experience

Time and location (Internationalization)



Weather (weather.com)

Personalized user experience

Past user search queries



User privacy

Backend development

- Data collection
- Data storage
- Configuration
- Logging
- Dashboard
- Security

Data collection

Data ownership

- Internal data (e.g., official website, internal wikis, databases etc.)
- External data (e.g., other websites, wikis, open data)

Data model (Data and Schema)

- Unstructured data (e.g., documents, texts, web pages etc.)
- Semi-structured data (e.g., JSON/XML files etc.)
- Structured data (e.g., relational databases, linked data)

Data collection

Data sources

- Web pages
- Documents, texts
- Sensors
- Databases
- ...

Data collection

Data acquisition

- Data dumps
- Crawlers
- Web scraping
- Application Programming Interface (API)

Data collection

Data cleaning and transformation

- Accuracy (e.g., verification with external sources)
- Validity (e.g., detect constraint violations)
- Uniformity (e.g., units)

Data storage

- Model
- Indexation
- Query optimization
- Caching
- Replication
- Backup

Data Model

- Database schema
- Schema-less

Data storage

- Relational Databases
- Object-oriented Databases
- NoSQL databases (e.g., graph databases)
- NewSQL databases (SQL + ACID guarantees)

Document indices and Query Optimization

Document indices

- Forward index
- Inverted index

- Join ordering
- Cost estimation

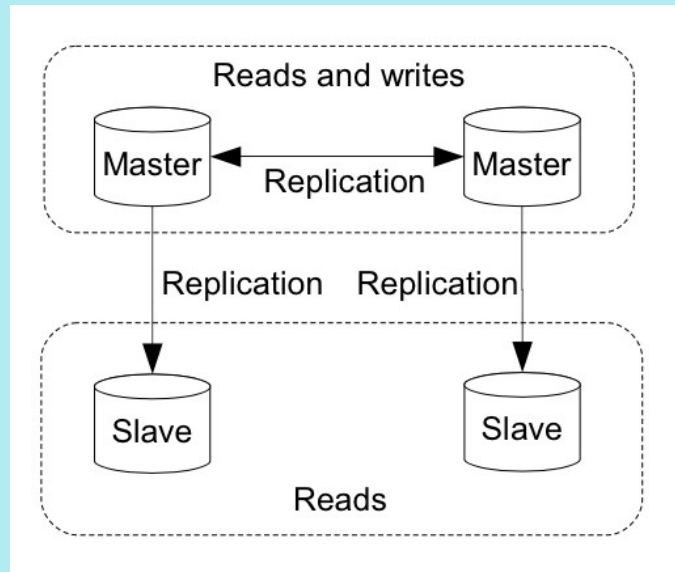
Database Indexation

Query Optimization

Caching

- Frequently asked questions and cached responses

Replication and Backup



Replication(Master-slave)

Resource management and configuration

Availability % ↕	Downtime per year ↕	Downtime per month ↕	Downtime per week ↕	Downtime per day ↕
90% ("one nine")	36.5 days	72 hours	16.8 hours	2.4 hours
95% ("one and a half nines")	18.25 days	36 hours	8.4 hours	1.2 hours
97%	10.96 days	21.6 hours	5.04 hours	43.2 minutes
98%	7.30 days	14.4 hours	3.36 hours	28.8 minutes
99% ("two nines")	3.65 days	7.20 hours	1.68 hours	14.4 minutes
99.5% ("two and a half nines")	1.83 days	3.60 hours	50.4 minutes	7.2 minutes
99.8%	17.52 hours	86.23 minutes	20.16 minutes	2.88 minutes
99.9% ("three nines")	8.76 hours	43.8 minutes	10.1 minutes	1.44 minutes
99.95% ("three and a half nines")	4.38 hours	21.56 minutes	5.04 minutes	43.2 seconds
99.99% ("four nines")	52.56 minutes	4.38 minutes	1.01 minutes	8.64 seconds
99.995% ("four and a half nines")	26.28 minutes	2.16 minutes	30.24 seconds	4.32 seconds
99.999% ("five nines")	5.26 minutes	25.9 seconds	6.05 seconds	864.3 milliseconds
99.9999% ("six nines")	31.5 seconds	2.59 seconds	604.8 milliseconds	86.4 milliseconds
99.99999% ("seven nines")	3.15 seconds	262.97 milliseconds	60.48 milliseconds	8.64 milliseconds
99.999999% ("eight nines")	315.569 milliseconds	26.297 milliseconds	6.048 milliseconds	0.864 milliseconds
99.9999999% ("nine nines")	31.5569 milliseconds	2.6297 milliseconds	0.6048 milliseconds	0.0864 milliseconds

Availability (Wikipedia)

Resource management and configuration

- Machines (servers, disks etc.)
- Software packages and dependencies
- Energy consumption

Deployment

- Development setup
- Pre-production setup
- Production setup

Packaging

- Containers (e.g., Linux containers)

Load balancing

- Server-side
- Client-side

Selective Testing

A/B Testing

Logging

- Access logs
- Error logs
- Event logs
- Transaction logs

Logging

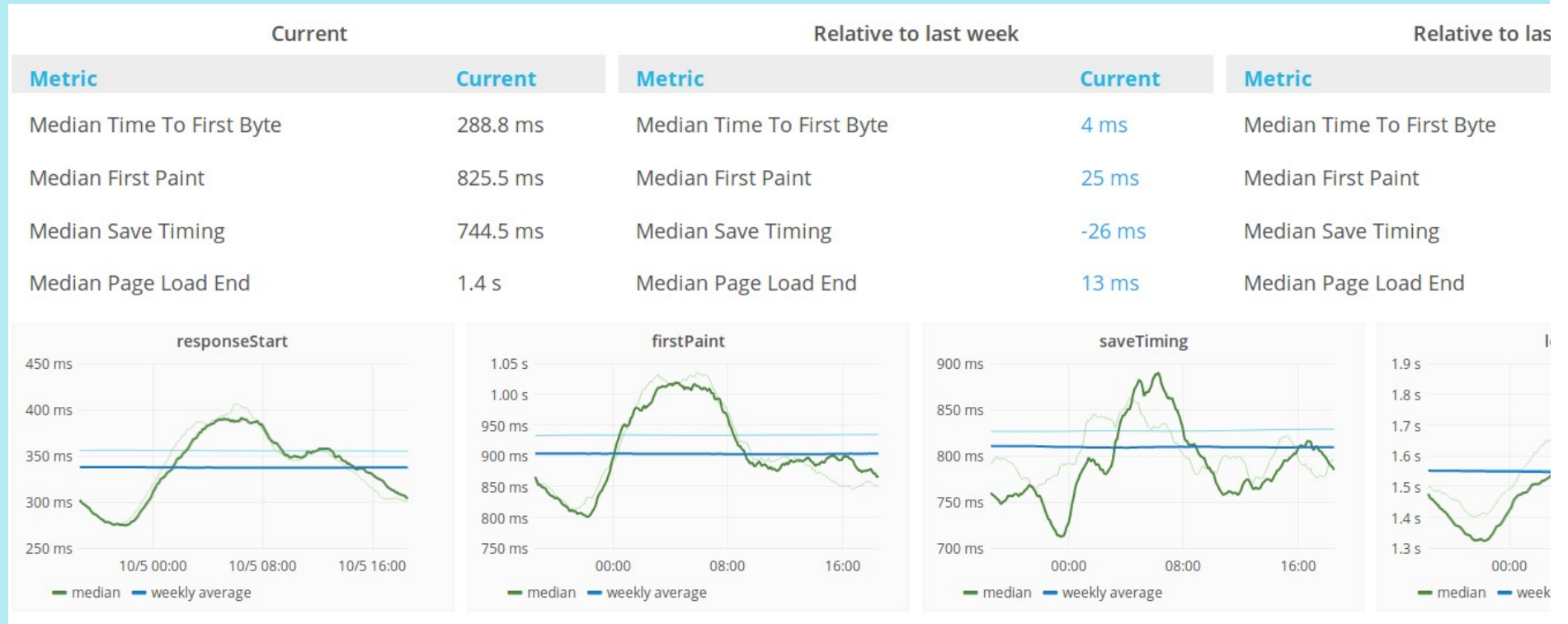
Why logs?

- Debugging
- Security (e.g., detect intrusion)
- Database rollbacks
- Audit
- Analysis (e.g., detecting patterns, resource planning)

Logging

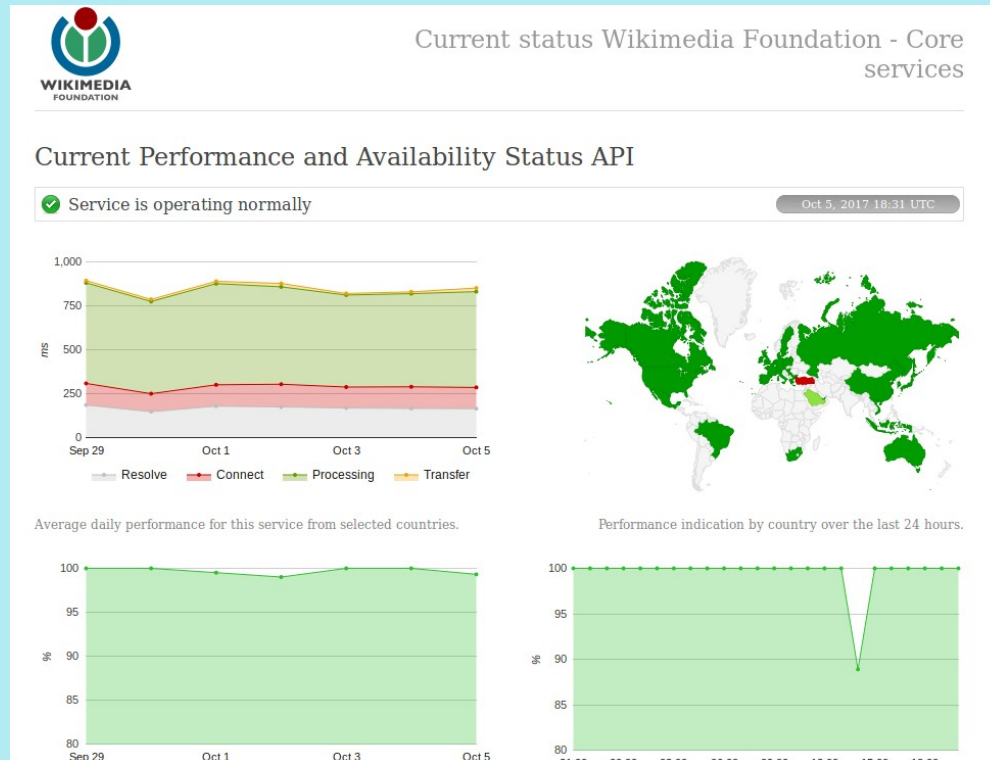
- IP address
- User
- Resource ID
- ...

Dashboard



Wikimedia (Grafana: 5th October 2017)

Dashboard



Wikimedia (Availability: 5th October 2017)

Dashboard

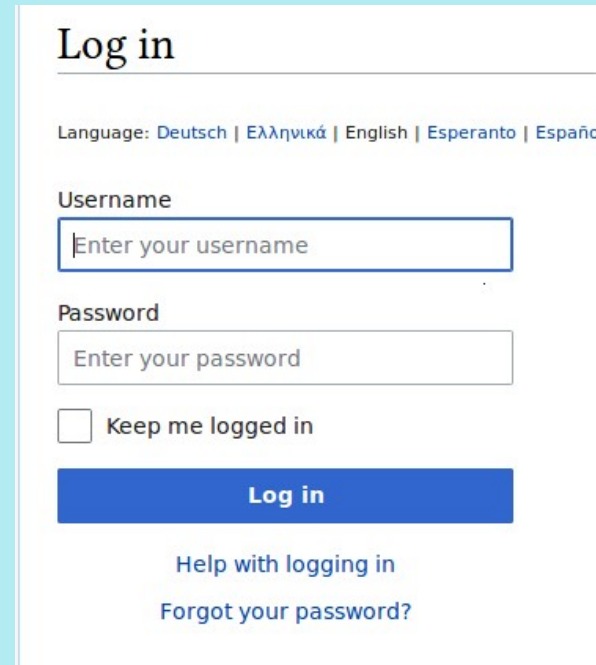
- Performance Metrics (active users, queries served etc.)
- Real-time metrics (e.g., downtime, latency, throughput)

Dashboard

- Email Alerts
- Visual indicators

Security

- Data protection
- Logged-in users or public access
- Third party access

A screenshot of the Wikipedia login page. At the top, it says "Log in" in a large, dark font. Below this, there is a language selection bar with links for "Deutsch", "Ελληνικά", "English", "Esperanto", and "Español". The main form has two input fields: "Username" with a placeholder "Enter your username" and "Password" with a placeholder "Enter your password". Below the password field is a checkbox labeled "Keep me logged in". A blue "Log in" button is positioned below the checkbox. At the bottom of the form, there are two links: "Help with logging in" and "Forgot your password?".

Log in

Language: [Deutsch](#) | [Ελληνικά](#) | [English](#) | [Esperanto](#) | [Español](#) | [Français](#) | [Galego](#) | [Italiano](#) | [日本語](#) | [Polski](#) | [Português](#) | [Română](#) | [Slovenščina](#) | [Tagalog](#) | [Türkçe](#) | [Українська](#) | [Vietnamese](#) | [Walon](#) | [Yorùbá](#)

Username

Enter your username

Password

Enter your password

☐ Keep me logged in

Log in

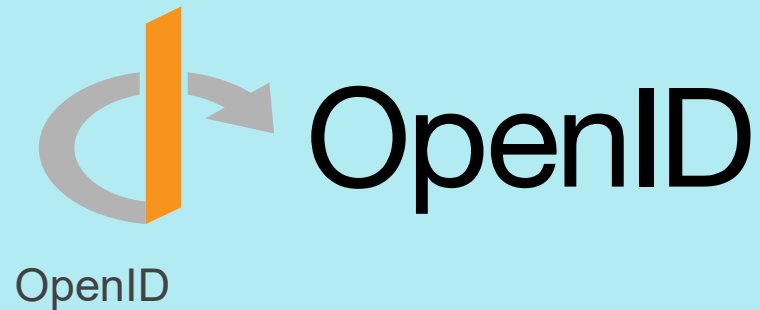
[Help with logging in](#)

[Forgot your password?](#)

Login (Wikipedia)

Security

- Authentication
- Authorization to third party access



Mozilla Persona	
	mozilla Persona
Developer(s)	Mozilla Foundation
Initial release	July 2011
Written in	JavaScript
Operating system	Cross-platform
Available in	51 languages
Type	Authorization
License	MPL
Website	developer.mozilla.org/en-US/Persona

Mozilla Persona (2011-2016)

Detecting security vulnerabilities

- Intrusion
- SQL code injection
- Cross-site scripting
- Denial of service

Application programming interface

- Service-oriented (SOAP)
- Resource-oriented (REST)

API: Data formats

- XML
- JSON

API: (CRUDL) Operations

- Create
- Read
- Update
- Delete
- List

API: Data dumps

- Complete data dumps
- Selective data dumps

Application programming interface

- HTTP
- Software development kits (SDKs)

Interface definition

- Human-readable documentation
- Machine-readable documentation (WSDL, WADL etc.)
- Human and machine-readable documentation (microformats, semantic web languages)

Human readable Documentation

1. Read documentation
2. Develop application to integrate
3. Add business logic, if any

Machine-readable Documentation

1. Fully autonomous solution to integrate
2. Add business logic, if any

Quality of service

Resource usage limits

- Limits on API call count (per user, IP)
- Limits on data transfer
- Temporary blocks

Quality of service

- Analysis on frequently made API calls
- Resource planning and allocation

Security

- No password
- Basic authentication (e.g., username, password)
- (Open) authentication protocols (e.g., OAuth)



OAuth

Project



Virtual Library

Project



Target audience

Project

Search

Search

Time

Location

...

References

References

- https://en.wikipedia.org/wiki/Information_system
- [https://en.wikipedia.org/wiki/Search_engine_\(computing\)](https://en.wikipedia.org/wiki/Search_engine_(computing))
- <https://query.wikidata.org/>
- <https://weather.com/>
- <https://duckduckgo.com/bang>
- <http://commoncrawl.org/>
- https://en.wikipedia.org/wiki/High_availability
- <https://grafana.wikimedia.org>
- <https://status.wikimedia.org>
- <http://highscalability.com/>

Image credits

- [Wikimedia Commons](#)
- <http://histropedia.com/timeline/5b>
- <https://pixabay.com/>