# IANCIS
## Indexing of Anonymous Networks for Crime Information Search

Alessandro Celestini
Institute for Applied Computing (IAC-CNR)
a.celestini@iac.cnr.it

May 20, 2015

# IANCIS

**Objective**:   Developing a tool able to crawl onion websites and feed a semantic engine for indexing and clustering collected data.

**Website**:  www.iancis.eu

**Consortium**:

Istituto per le Applicazioni del Calcolo "Mauro Picone"

Arma dei Carabinieri HQs-ICT Office

Expert System

# IANCIS

**The tool should:**

- Automatically **explore** TOR hidden services

- **Extract** text from collected resources and **analyze** it

- **Visualize** analysis results and **identify** evidences of illegal activities

*The tool is meant to be the first step towards the development of an investigation instrument for the TOR network*
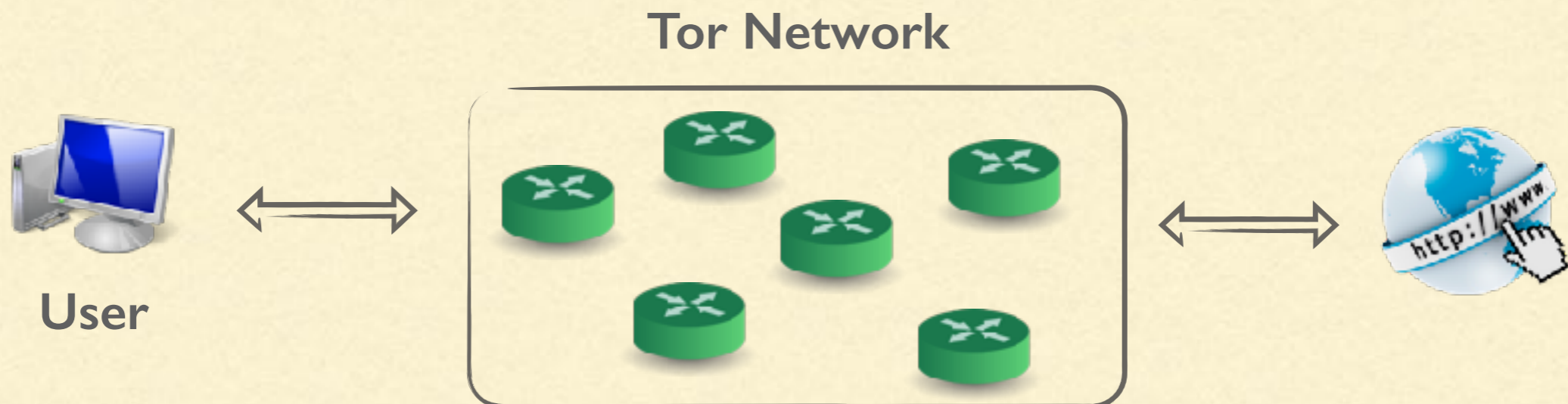
# TOR

Tor is an **anonymity network**, it was originally designed, implemented, and deployed as a third-generation **onion routing** project of the U.S. Naval Research Laboratory.

**Tor network** is composed by volunteer-operated servers (relays) that allow people to improve their privacy and security on the Internet.

Tor protects users against a common form of Internet surveillance known as **"traffic analysis"**. Traffic analysis can be used to infer **who is talking to whom** over a public network.

# TOR

Tor is an **anonymity network**, it was originally designed, implemented, and deployed as a third-generation **onion routing** project of the U.S. Naval Research Laboratory.

**Tor network** is composed by volunteer-operated servers (relays) that allow people to improve their privacy and security on the Internet.

Tor protects users against a common form of Internet surveillance known as **"traffic analysis"**. Traffic analysis can be used to infer **who is talking to whom** over a public network.

**Tor Network**



**User**

# Hidden Services

Tor can provide **receiver privacy** for Internet services through a feature called "hidden services".

Tor's hidden services let users publish web sites and other services **without revealing the location of the site.**
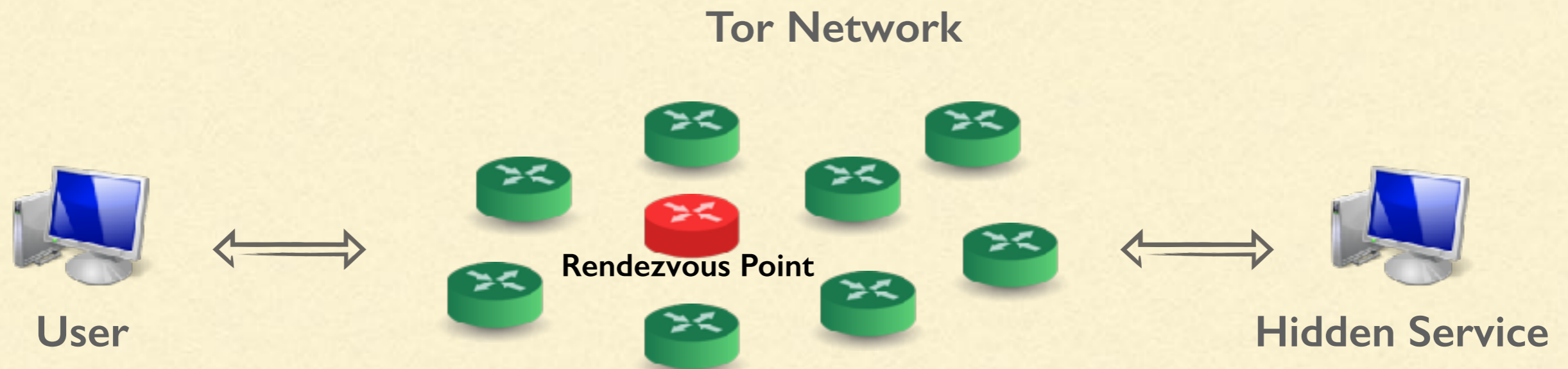
# Hidden Services

Tor can provide **receiver privacy** for Internet services through a feature called "hidden services".

Tor's hidden services let users publish web sites and other services **without revealing the location of the site.**

Hidden services are accessible only through the Tor network and each service is identified by its *onion address.* An onion address is a **16 characters** string followed by the domain .onion, e.g http://duskgytldkxiuqc6.onion.

The number of Hidden Services is estimated to be in the order of **tens of thousands**[1,2]

[1] Torproject Blog, "Some statistics about onions", https://blog.torproject.org/blog/some-statistics-about-onions
[2] A. Biryukov, I. Pustogarov, and R. Weinmann. 2013. Trawling for Tor Hidden Services: Detection, Measurement, Deanonymization. In Proceedings of SP '13

# Hidden Services

Tor can provide **receiver privacy** for Internet services through a feature called "hidden services".

Tor's hidden services let users publish web sites and other services **without revealing the location of the site.**

Hidden services are accessible only through the Tor network and each service is identified by its *onion address.* An onion address is a **16 characters** string followed by the domain .onion, e.g http://duskgytldkxiuqc6.onion.

The number of Hidden Services is estimated to be in the order of **tens of thousands**[1,2]

## Tor Network

**Rendezvous Point**

**User**

**Hidden Service**

[1] Torproject Blog, "Some statistics about onions", https://blog.torproject.org/blog/some-statistics-about-onions
[2] A. Biryukov, I. Pustogarov, and R. Weinmann. 2013. Trawling for Tor Hidden Services: Detection, Measurement, Deanonymization. In Proceedings of SP '13

# Hidden Services

**Why LEA are interested in Hidden Services?**

Several Hidden Services **promote illegal activities** such as:

- Illicit trafficking in weapons
- Illicit trafficking in narcotic drugs
- Child pornography
- Trafficking in human beings
- Terrorism

# Hidden Services

**Why LEA are interested in Hidden Services?**

Several Hidden Services **promote illegal activities** such as:

- Illicit trafficking in weapons
- Illicit trafficking in narcotic drugs
- Child pornography
- Trafficking in human beings
- Terrorism

**Silk Road** was an online marketplace mostly used for selling illegal goods.
In 2013 the FBI shut down the website.

# Hidden Services

**Why LEA are interested in Hidden Services?**

Several Hidden Services **promote illegal activities** such as:

- Illicit trafficking in weapons
- Illicit trafficking in narcotic drugs
- Child pornography
- Trafficking in human beings
- Terrorism

**Silk Road** was an online marketplace mostly used for selling illegal goods. In 2013 the FBI shut down the website.

**Lolita City** hosted child pornography images and videos.

# Tool

# Workflow

Acquisition → Elaboration → Identification of Illegal activities

Elaboration → Visualization

# Workflow

# Workflow

# Workflow

# Workflow

Acquisition → Elaboration → Identification of Illegal activities

Elaboration → Visualization

# Technologies

ACQUISITION

ELABORATION **Text**

ELABORATION **Text**

# Technologies

ACQUISITION

**BUBING**: high performance crawler developed by the Laboratory for Web Algorithmics (LAW) of the University of Milan

ELABORATION **Text**

ELABORATION **Text**

# Technologies

ACQUISITION



**BUBING**: high performance crawler developed by the Laboratory for Web Algorithmics (LAW) of the University of Milan

ELABORATION **Text**



**TIKA**: open-source software suite for the identification and extraction of text from more than 1000 different file types

ELABORATION **Text**

# Technologies

ACQUISITION

**BUBING**: high performance crawler developed by the Laboratory for Web Algorithmics (LAW) of the University of Milan

ELABORATION **Text**

**TIKA**: open-source software suite for the identification and extraction of text from more than 1000 different file types

ELABORATION **Text**

COGITO

**COGITO**: semantic analysis engine by Expert System, that classifies a text according to a suitable taxonomy, providing both quantitative and qualitative information (*what topic* the text is about and *how much* the text discusses such topic)

# Architecture

Crawling

Semantic
Engine

Extraction and Filtering

Analysis and Filtering

**Extraction, Analysis and Storage**

WARC
(web archive)

Document
Oriented DB

Document
Oriented DB

# Architecture



Crawling

Semantic
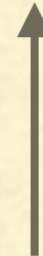Engine

Extraction and Filtering
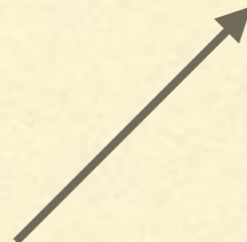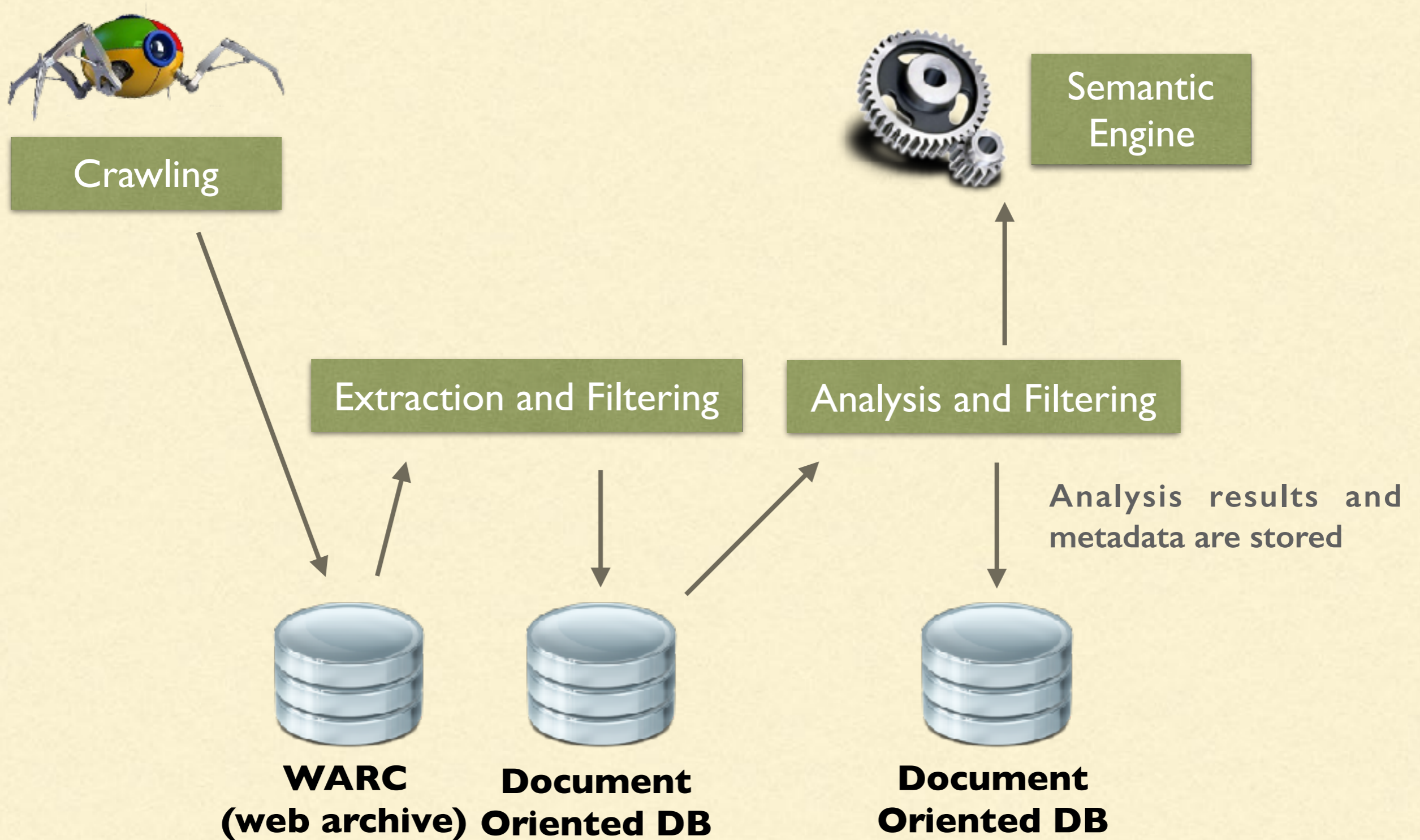
Analysis and Filtering

**WARC
(web archive)**

**Document
Oriented DB**

**Document
Oriented DB**

# Architecture

Crawling

Semantic
Engine

Crawled resources
are filtered and
stored for later
analysis

Extraction and Filtering

Analysis and Filtering

**WARC
(web archive)**

**Document
Oriented DB**

**Document
Oriented DB**

# Architecture



Crawling

Semantic Engine

Extraction and Filtering

Analysis and Filtering

**WARC (web archive)**

**Document Oriented DB**

**Document Oriented DB**

# Architecture

**Crawling**

**Semantic Engine**

Text is extracted and sent for analysis

**Extraction and Filtering**

**Analysis and Filtering**

**WARC (web archive)**

**Document Oriented DB**

**Document Oriented DB**

# Architecture



Crawling

Semantic Engine

Extraction and Filtering

Analysis and Filtering

**WARC (web archive)**

**Document Oriented DB**

**Document Oriented DB**

# Architecture



Crawling

Semantic Engine

Extraction and Filtering

Analysis and Filtering

Analysis results and metadata are stored

**WARC (web archive)**

**Document Oriented DB**

**Document Oriented DB**

# Architecture

# Visualization

# Tag Cloud

Shows the **entities extracted** by Cogito
(in the documents retrieved by the crawler)



Three types of entities are identified by the engine ⟶  **People**
**Places**
**Organizations**

# Tag Cloud

Shows the **entities extracted** by Cogito
(in the documents retrieved by the crawler)



Three types of entities are identified by the engine → **People Places Organizations**

# Tag Cloud

Shows the **entities extracted** by Cogito
(in the documents retrieved by the crawler)



Three types of entities are identified by the engine  ⟶  **People**
**Places**
**Organizations**

# Treemap

Shows **documents by category**
(categories are assigned by the engine)

# Treemap

Shows **documents by category**
(categories are assigned by the engine)

# Treemap

Shows **documents by category**
(categories are assigned by the engine)

# Thank you for your time