# FANDANGO Project: Advanced analytics services to detect disinformation

**ZAGABRIA, NEM SUMMIT 2019 , 23TH MAY 2019**

# How to spot fake-news? According to IFLA

**THE INTERNATIONAL FEDERATION OF LIBRARY ASSOCIATIONS AND INSTITUTION**

*(based on FactCheck.org's 2016 article How to Spot Fake News)*



## HOW TO SPOT FAKE NEWS

**CONSIDER THE SOURCE**
Click away from the story to investigate the site, its mission and its contact info.

**READ BEYOND**
Headlines can be outrageous in an effort to get clicks. What's the whole story?

**CHECK THE AUTHOR**
Do a quick search on the author. Are they credible? Are they real?

**SUPPORTING SOURCES?**
Click on those links. Determine if the info given actually supports the story.

**CHECK THE DATE**
Reposting old news stories doesn't mean they're relevant to current events.

**IS IT A JOKE?**
If it is too outlandish, it might be satire. Research the site and author to be sure.

**CHECK YOUR BIASES**
Consider if your own beliefs could affect your judgement.
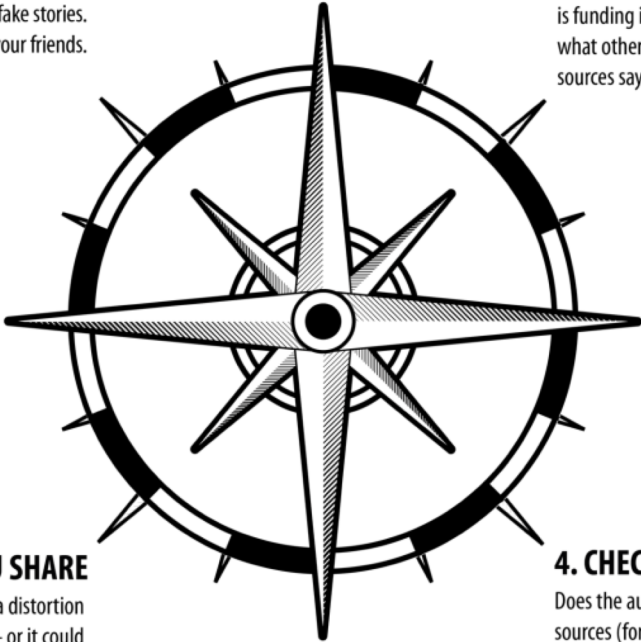
**ASK THE EXPERTS**
Ask a librarian, or consult a fact-checking site.

IFLA
International Federation of Library Associations and Institutions

# How to spot fake-news?
# According to EPRS

**EUROPEAN PARLIAMENTARY RESEARCH SERVICE**



**1. CHECK THE CONTENT**
Are the facts and figures accurate? Is the article biased? A credible media outlet keeps one-sided opinions where they belong – in op-eds, not in news articles.

**2. CHECK THE OUTLET**
Do you know it? Does the URL look strange? Check the 'about' section. Who is behind it? Who is funding it? Double-check what other (trustworthy) sources say.

**3. CHECK THE AUTHOR**
Does this person even exist? A well-respected journalist always has a track record. If the author has made up his or her name (or does not mention it), the rest is also likely to be fake.

**4. CHECK THE SOURCES**
Does the author use reliable sources (for example, well-established and respected media outlets)? Are the quoted experts real specialists? If the story uses anonymous (or no) sources, it could be fake.

**5. CHECK THE PICTURES**
Images are powerful, and it is easy to manipulate them. An image search can show if it has been used before in a different context. The InVID plugin[1] can help you detect manipulation of videos or pictures.

**6. THINK BEFORE YOU SHARE**
The story could be a distortion of real or old events – or it could be satire. The headline could be designed to spark strong emotions. If an event is real, reliable media will cover it.

**7. QUESTION YOUR OWN BIASES**
Sometimes a story is just too good or entertaining to be true. Take a deep breath, compare with reliable sources and keep a cool head.

**8. JOIN THE MYTH-BUSTERS[2]**
Keep on top of the latest tricks and narratives used by those spreading disinformation. Report fake stories. Tell your friends.

ENGINEERING | LIVE Tech LIVING TECHNOLOGY | CERTH CENTRE FOR RESEARCH & TECHNOLOGY HELLAS | SIREN DATA INTELLIGENCE | vrt | CIVIO | POLITÉCNICA | ANSA

# How Fandango aims to tackle disinformation?

providing an online service that will support professionals with the following features:

▶ News disinformation detection and scoring, based on Big Data analysis techniques (ML models and Graph Analysis)

▶ data investigation, through an interactive exploration of news, open data and verified claims databases.

# Disinformation scoring features

Fandango provides a set of disinformation scoring features by analyzing the different components of news:

▶ Text (headline, body)

▶ Authors & Source

▶ Media (images, videos)

# Disinformation scoring features
# Text analysis: our approach (ML)

► machine learning model will be trained to recognize features in the new's headline and body

► We are testing a completely "**context agnostic**" classifier, it means that words are not considered as single feature in the model.

► We are testing different machine learning algorithms in a greedy way, with an accurate features selection. It means that we train lot of models, from simple to complex ones, an take the one with best performs in terms of accuracy and precision.

# Disinformation scoring features
# Text analysis: our ML features (1/2)

Features **applied** to headline and body of article:

**Simple frequency features**

▶ Counting Stopwords - the stopwords are the most frequent word in a given language
(e.s the, and, or, with.. ciao, il, la, ... ola etc.)

▶ Characters Counter

▶ Punctuation Counter

**Part of The Speech (POS) features:**

▶ Counting POS( adjectives, adverbs, verbs, conjunctions)

# Disinformation scoring features
# Text analysis: our ML features (2/2)

Features **applied** to title and text of an article are:

**Advance frequency features**

- ▶ Word average per paragraph
- ▶ Lexical Diversity

**readability indices:** readability tests designed to assess the understandability of a text.

- ▶ Flesch Index of readability
- ▶ FKG Read Level
- ▶ Automated Readability Index(ARI)
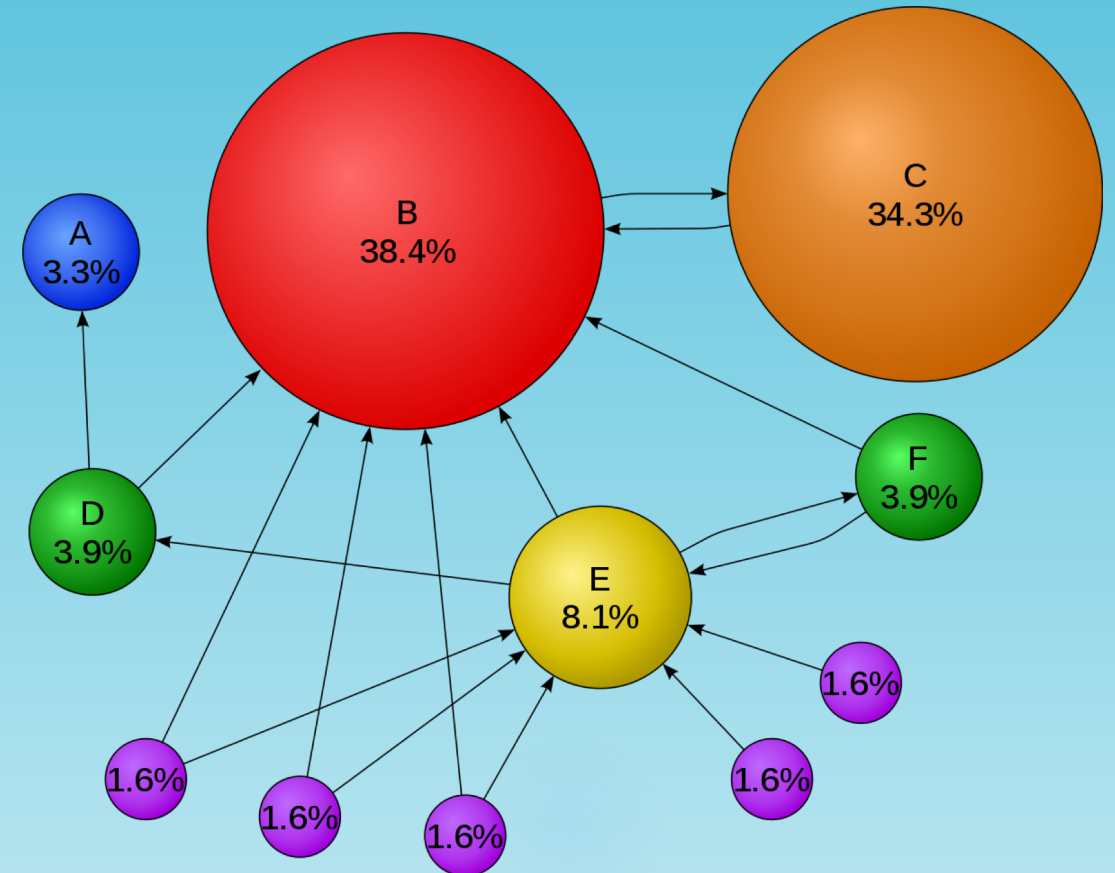
# Disinformation scoring features
# Text analysis: differences in term of features

# Disinformation scoring features
# Authors & Source analysis: our approach (GA)

► We apply importance algorithms to get the impact of each entity (e), including authors and organizations, in FANDANGO's network.

► Weighted credibility indicator using the disinformation score of the articles provided by the text analyzer.

► Scoring result: will be between 0 and 1 if we have information enough for a particular entity, otherwise it will yield -99 to remark the lack of information to compute the analysis.

# Disinformation scoring features
# Authors & Source analysis: our approach (GA)

► In one hand, supervised learning has the most impressing results in the deep learning methodologies

► On the other hand, the problem with big data is that it's impossible to have fully annotated datasets to train the DL models in that manner

► a semi-supervised approach is the most realistic in terms of annotation effort / time compared to the final outcomes

► Graph based techniques are very effective in representing linearly non-separable data, but there is a limitation in working in a full adjacency matrix of the data, which translates to some millions of elements in a large dataset

► to overcome this issue we examined a sampling / batch learning approach that works quite effectively both in small scale and (most importantly) in large scale graphs

# Disinformation scoring features

## Image & Video analysis

### Spatiotemporal Analytics and out of context

Our goal is the detection of out of context content.
Our approach is based on comparing topics and entities extracted from the body of an article with the topic extracted from an image or video.

we are focused to detect 3 main types of out-of-context:
- two news about different topics but containing the same image or a manipulated version (i.e. using the same guerrilla images in two different war scenarios)
- two different news share the same image (or a manipulated version) but the publication dates are far from each other
- two different news share the same image (or a manipulated version) but the news are about two different location entities.

# Disinformation scoring features
# Image & Video analysis

## Spatiotemporal Analytics and out of context

The technique is to apply LDA NLP (Latent Dirichlet allocation - natural language processing) methods. This is the most common topic model currently in use, based on the intuition that documents cover a small number of topics and that topics often use a small number of words.

**Input:** News articles (**DATASET**) gathered from the web (as json files) and stored in 'Articles DB'.

**Output:** Entities, leading words, topics (for each article)

# Disinformation scoring features
# Image & Video analysis

## Spatiotemporal Analytics and out of context: Object Detection

Use pre-trained model to detect objects

Create a vector with all detected objects (probability greater than 30%)

This vector describes the content of the image

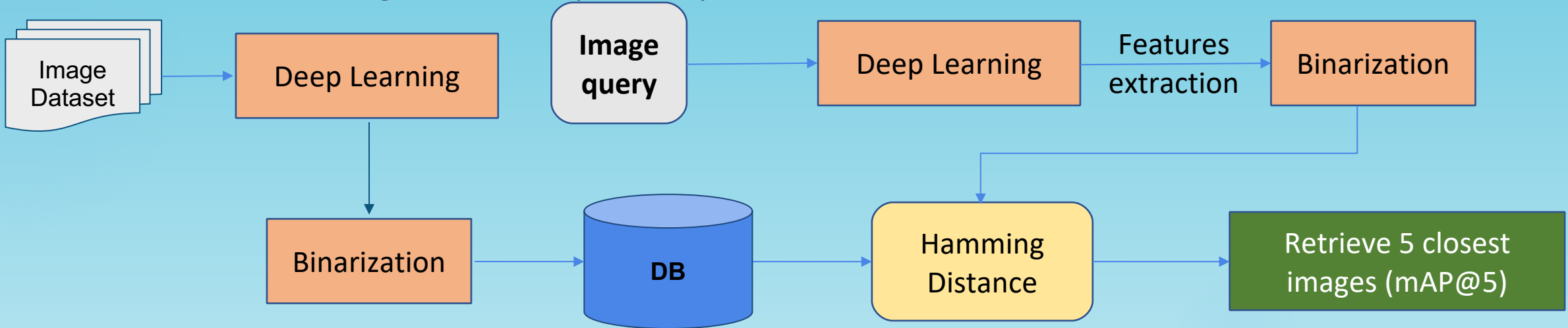# Disinformation scoring features
# Image & Video analysis

## Spatiotemporal Analytics and out of context: Image Similarity

Test Datasets = synthetic, ukbench, copydays, mscoco
Queries =150
Image database = 10k
Mean average precision (mAP@5)

# Results Visualization 1/2

❖ **MSCOCO**
**gaussian noise**

❖ **COPYDAYS**
**50% crop**

❖ **MSCOCO**
**vertical flip**



Query Results [3, 7, 9, 9, 9]

Query Results [6, 7, 8, 8, 8]

Query Results [8, 9, 9, 9, 10]

# Results Visualization 2/2

❖ **MSCOCO horizontal flip**


Query Results [3, 5, 8, 8, 8]

❖ **UKBENCH Dataset**


Query Results [2, 5, 5, 5, 5]

# Disinformation scoring features
# Image & Video analysis
## Copy-Move Detection

**Goal.** Build a model to classify images as fake or pristine. Provide visualizations showing the fake objects

**Model.** CNN Binary Classifier

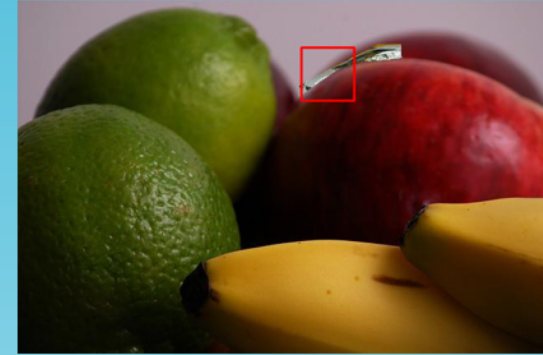**Datasets.** CERTH synthetic (~2K), CASIA small (~2K), CoMoFoD (200), 1st Image Forensics Challenge (~2K)
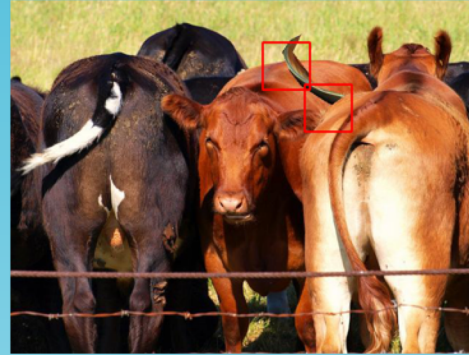
**Preprocessing.** Patch extraction using binary masks for random sampling on pristine, sampling on the boundaries of fake objects for fake images.

**Evaluation.** Perform the exact same preprocessing to test images to obtain fake and pristine patches. Evaluate images as fake even if a single patch is predicted as fake with probability > threshold. Pristine if no patch is predicted as fake.

# Copy-Move Detection:Visualizations

*(Fakeness) probability threshold: 0.8. Patches that got high probability are marked with a red square*



Evaluation on Fake Images

# Copy-Move Detection:Visualizations

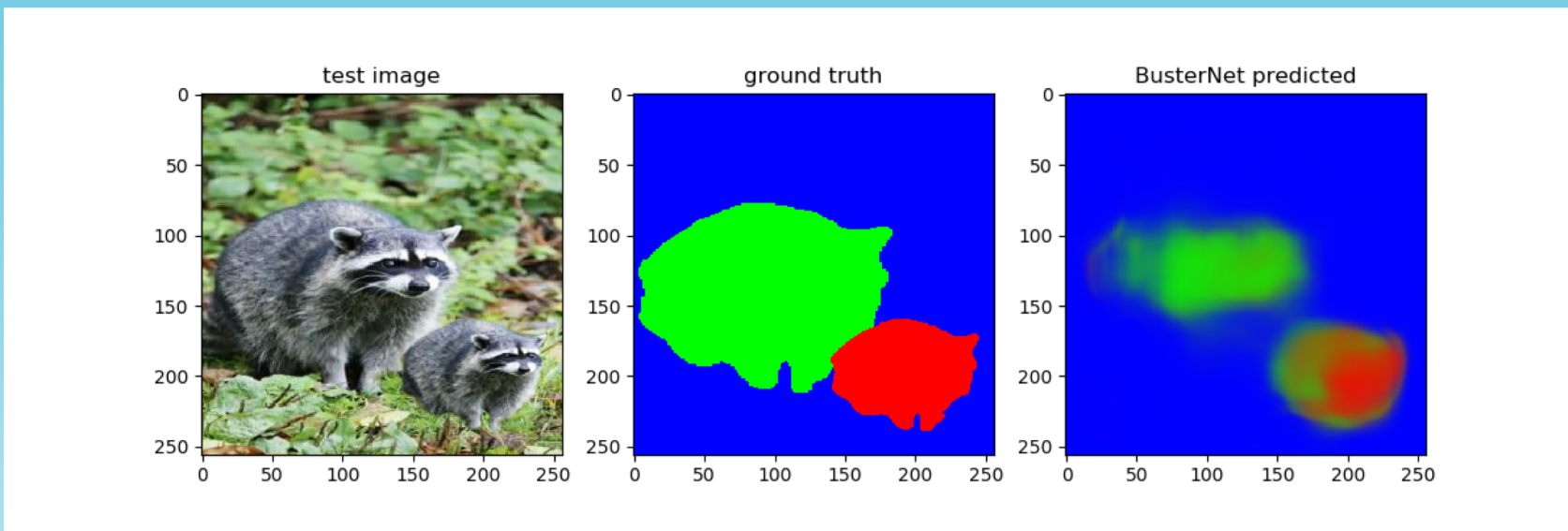*(Fakeness) probability threshold: 0.8. Patches that got high probability are marked with a red square*

Evaluation on Pristine Images

# Copy-Move Detection: Visualizations

The final tool that will be obtained, is a multi – class classifier, able to predict three – class masks for the given fake, or pristine images

# Data Investigation features
# Claim analysis: claim reviews

GOAL: support data investigation, providing similar claim reviews of a examined claim

FANDANGO will search its internal Claim database, collected from trustworthy sources that provide Claim Reviews, and display the most similar Claims and its associated Claim Reviews

It performs text comparison and tf/idf similarity analysis, which provides solid results on western languages, combined with a custom weighted pipeline to provide an overall similarity score of the Claim.

# Data Investigation features
# Claim analysis: Open Data references

GOAL: support data investigation, providing direct links to Open Data
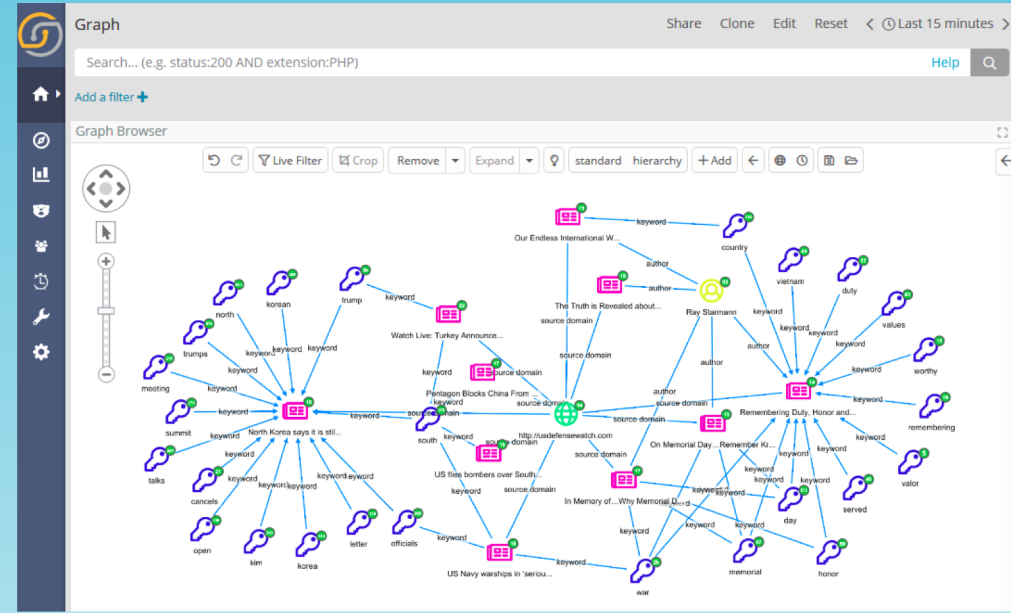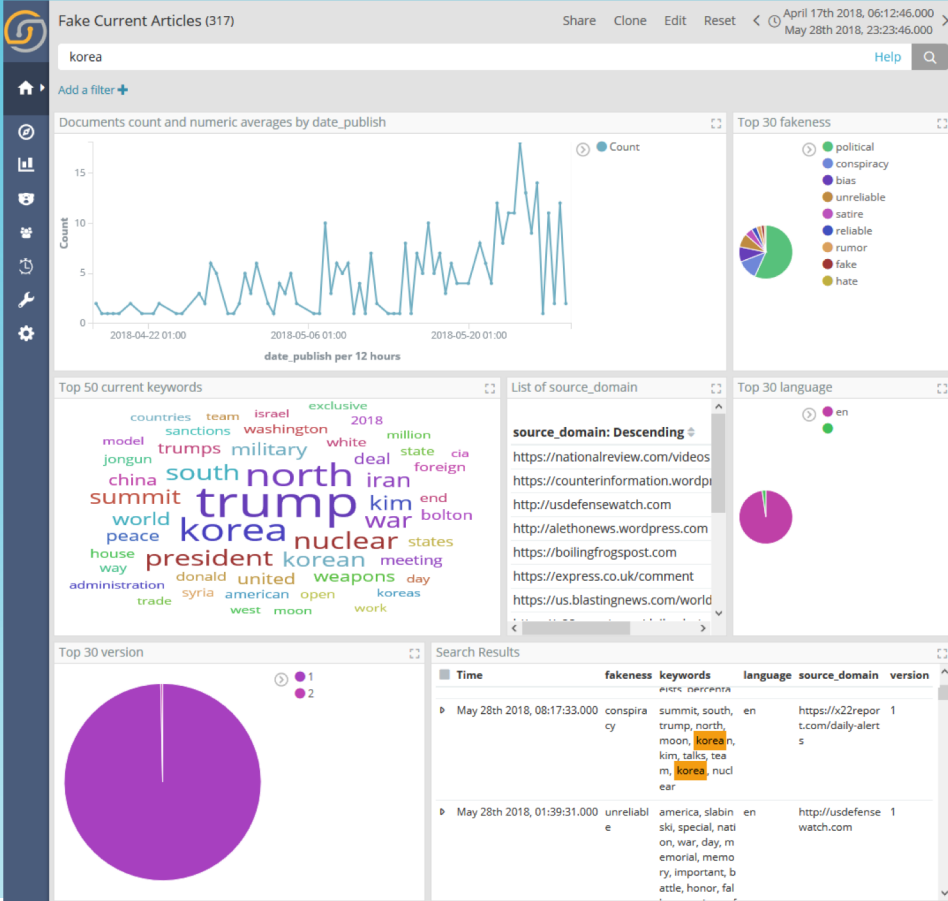
A list of links to Open Data references, with its titles, will be provided to the user based on a Claim selection or Article topic.

It performs a text comparison and tf/idf similarity analysis combined with a custom weighted pipeline to provide an overall similarity score of the Claim and topics.
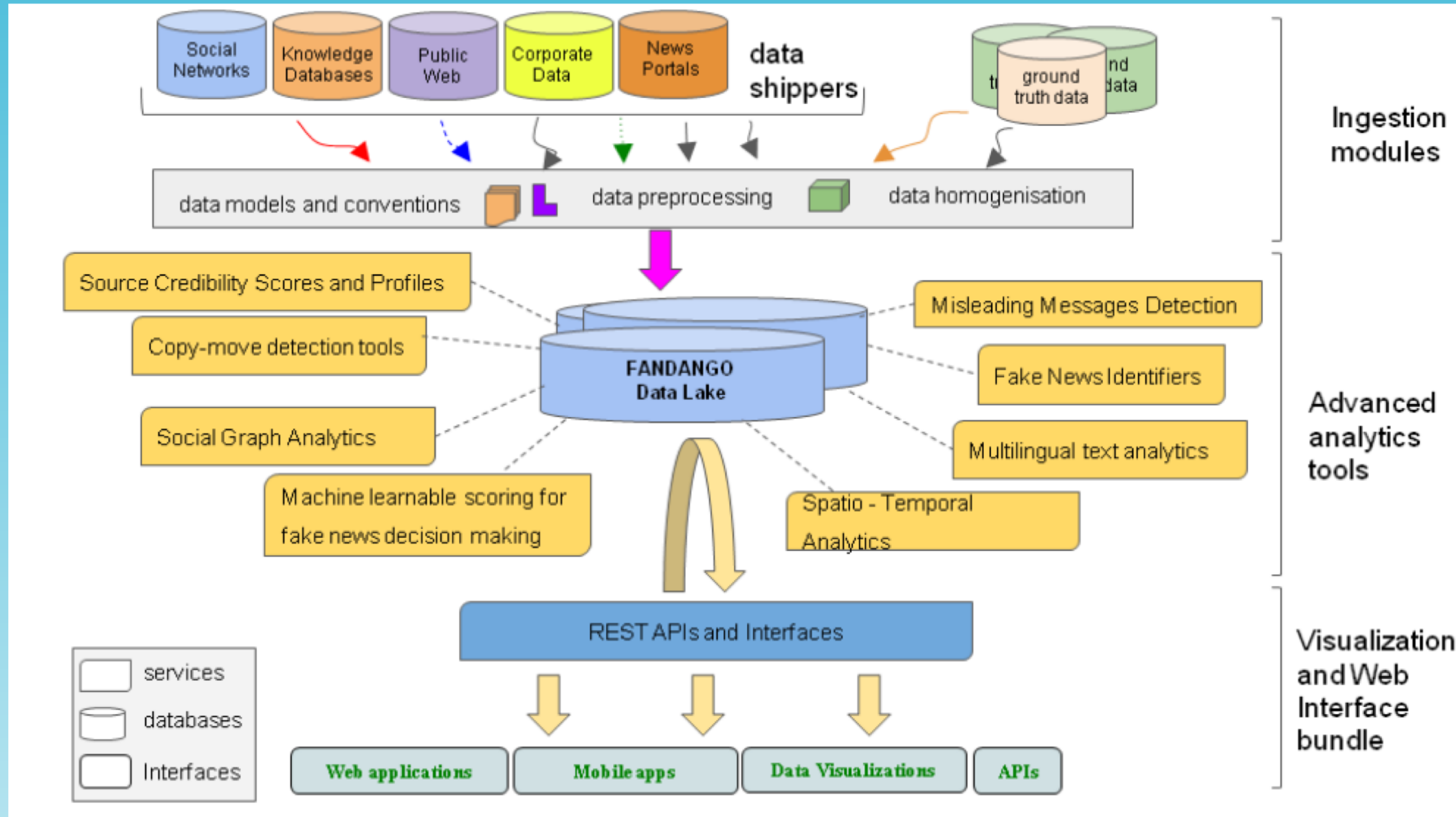
# Data Investigation features
# Dashboards & Knowledge Graph

# Architectural approach

# Developing towards FANDANGO End Users

- **Checking Against Previous Fact Checks**: FANDANGO will match statements to previous fact-checks or consulting authoritative sources

- FANDANGO will parse statements in terms that **make sense** to a database.

- FANDANGO will help in the **verification process** itself, to find the right data, the references, assess sources/publishers, etc.

- FANDANGO is focusing on a **user friendly set of tools** to assist journalist and will leave the **final decision** about the trustworthiness of information to the journalists.

Twitter: @fandango_eu
Facebook: @fandango.project

massimo.magaldi@eng.it


Thank you!