# K-POP and Fake Facts: from Texts to Smart Alerting for Maritime Security

**Maxime Prieur**[1,2] and **Souhir Gahbiche**[1] and **Guillaume Gadek**[1]
**Sylvain Gatepaille**[1] and **Killian Vasnier**[1] and **Valerian Justine**[1]

(1) Airbus Defence and Space, 1 Bd Jean Moulin, 78990 Elancourt, France
(2) CNAM, 292 rue Saint-Martin, 75003 Paris, France
`[maxime.prieur, souhir.gahbiche, guillaume.gadek]@airbus.com`

## Abstract

Maritime security requires full-time monitoring of the situation, mainly based on technical data such as radar or Automatic Identification System (AIS) but also from Open Source Intelligence like inputs (e.g., newspapers). Some threats to the operational reliability of this maritime surveillance, such as malicious actors, introduce discrepancies between hard and soft data (sensors & texts), either by tweaking their AIS emitters or by emitting false information on pseudo-newspapers.

Many techniques exist to identify these pieces of false information, including using knowledge base population techniques to build a structured view of the information. This paper presents a use case for suspect data identification in a maritime setting. The proposed system UMBAR ingests data from sensors and texts, processing them through an information extraction step, in order to feed a Knowledge Base (KB) and finally perform coherence checks between the extracted facts.

## 1 Introduction

One of the main challenges in the maritime domain is to ensure safety and security of ships: in and around harbors but also when they are offshore for several days. The security aspect has benefited from a renewed interest recently, due to piracy and trafficking. Most harbor administrations rely today on AI-powered investigation tools to perform a number of checks on each entering vessel: comparing the declared status of the ship, aggregating sensor data, and even searching the Web for news.

Among the organizations that collect and disseminate information about maritime events, the Maritime Information Cooperation and Awareness Center (MICA) collects and relays useful information to all actors in the field of maritime industry. Its purpose is to process maritime security data worldwide. The 2022 annual report[1] summarizes the reports regularly sent to the maritime industry and analyses the trends observed as well as the evolution of modes of action.

The sensor data mainly consists of radar and AIS signals. Every ship must emit its identity, speed, position and course at short time intervals. This information is received by all other vessels in reach as well as dedicated receiving stations, on the coast and in space.

Based on sensor data, alerts related to the behavior of vessels are raised automatically: abnormal position, sudden change of direction, etc. This ensures a quick and efficient reaction of the police/security agents.

Relevant information about ships and maritime events also occurs in a non-technical way, through the news (so-called "*soft data*"). Accidents, illegal events, presence of a vessel in blockade-regulated areas and even modifications in the financial structure of the proprietary company are highly susceptible to increase the risk of a ship entering a harbor.

Malicious actors may use a large variety of techniques in order to perform covert operations, including trafficking, illegal fishing, piracy and smuggling. AIS are easy to tamper with, as ships may -illegally- decide to modify their identity, declare a false destination or even cease to emit.

Another case of concern occurs when civil vessels are the main object of a crisis between international powers, such as the *Stena Impero* near Iran in 2019 or the wheat vessels in the Black Sea in 2022: different newspapers may diffuse contradictory information about the same ships, which may in turn be inconsistent with technical data. In these cases, we believe that smart tools are needed in order to refine information and help the analysts build a clear picture of the situation.

To tackle these security challenges, we propose

---

[1] https://www.mica-center.org/en/home/download/2395/?tmstv=1673337653

UMBAR, a system to automatically collect, analyze and compare information from a variety of sources of data, resulting in a risk assessment that is practical for a security operator in the maritime domain. Such a system relies heavily on Natural Language Processing on the textual modality as well as on reasoning modules on the extracted knowledge.

More precisely, the contributions of this article are the following:

- a technical description of UMBAR, a complete operable system ranging from data collection to knowledge management,

- evaluation elements at a statistical and methodological level for the constituting subsystems,

- key points of attention towards a large-scale deployment of such a system.

The article is structured as follows: section 2 provides a review of the literature on the topic of AI-assisted maritime surveillance, with a special focus on knowledge based approaches; section 3 presents our system UMBAR and each of its subsystems from Information Extraction to Alert Raising; performance evaluation elements are provided at a subsystem level in section 4. A prospective discussion is exposed in section 5 to explicit the remaining challenges of the system deployment. Finally, section 6 concludes this paper.

## 2   Related Works

We structure our review of the literature along three streams: first, the identification of lies or manipulation on structured data. Detecting fake news has recently received a lot of attention, combining facts and language (Seddari et al., 2022); here we focus on the identification of dissimilarity between facts such as stored in knowledge bases (attributes, properties, relations). For media analysis purposes, this falls under the topic of "automatic fact-checking" (Guo et al., 2022).

In the maritime use case, expert systems for alert raising are common to detect a change of destination for a commercial vessel, or even a change of shipowner or flag can usually be observed (Alaeddine and Ray, 2022). AIS systems can also be hacked to disseminate false information manufactured. The objective of these false messages (e.g. distress signals, false vessel locations, etc.) is to attract attention and trap the targeted vessels (Balduzzi et al., 2014). These operations of disinformation and deception are very dangerous: it is essential to identify them.

Second, we focus our research on reasoning on facts in Knowledge Bases (KB), extracting the information using Knowledge base POPulation (K-POP), to automatically compute dissimilarity between text-extracted small Knowledge Graphs (KG) and to enable relation prediction; these applications are considered relevant for maritime security (Everwyn et al., 2019).

Zhang et al. (2019) focus on the link prediction task with complex linked datasets. Their approach successfully captures crossover interactions between entities and relations when modeling KGs. d'Amato et al. (2022) propose an approach based on semantic similarity for generating explanations to link prediction problems on Knowledge Graphs. Bhowmik and de Melo (2020) propose a model based on a Graph Transformer that learns entity embeddings by iteratively aggregating information from neighboring nodes to tackle the problem in the case of graphs that evolve over time.

Finally, our goal is to detect changes which occur over time and to evaluate information that evolves over time. Thus, this technique will identify a large number of alerts linked to a normal evolution of the characteristics of an entity. Dealing with temporal KB is still nothing trivial and mainly dealt with for Question-Answering where the relevant answers is dependent on time(Chen et al., 2022). Reasoning on such facts with intelligent systems is pretty much novel (Zhang et al., 2022), and still mainly dealt with by expert rules in operational systems.

## 3   System breakdown

In this section, we first sketch a view of UMBAR, then detail its two pillars: K-POP to extract the information, and coherence checks to perform the verification.

### 3.1   System Overview

The strength of our system, illustrated in the figure 1, lies in its ability to efficiently handle the end-to-end extraction and verification of valuable information from heterogeneous sources. The information contained in filtered sources is first extracted through the K-POP pipeline using transformer-based (Vaswani et al., 2017) language models. Extracted entities are compared to the existing ones
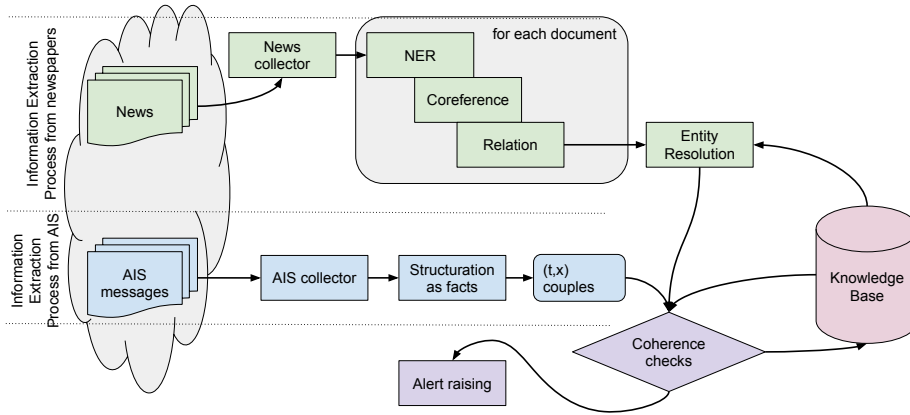
Figure 1: Overview of the UMBAR system: in green the information extraction pipeline from texts; in blue the processing of AIS messages; on the right a KB is fed as an output of the coherence check module.

in the KB to instantly detect and flag any inconsistencies. The end user is immediately aware of any potential alert and can then track down the cause from the relevant sources.

## 3.2 Information Extraction

Relevant information such as named entities (locations, organizations, persons and equipment), events and relations between these entities are extracted using a pipeline of Natural Language Processing (NLP) modules (Prieur. et al., 2023).

**Named Entity Recognition (NER):** The first component of this pipeline recognizes entities of interest in the text while assigning them a type with the help of the document. This block is instantiated by a fine-tuned RoBERTa (Liu et al., 2019) language model.

**Co-reference resolution:** It is then necessary to group the mentions referring to the same textual entities. In this case, the pre-trained World-level coreference resolution model (Dobrovolskii, 2021) is used to find groups of words referring to the same concept. These results are combined with those of the first block to obtain clusters with the same type.

**Relation extraction:** For this step we fine-tune the ATLOP model (Zhou et al., 2021) that produces an embedding of each entity at document scale before predicting the potential links (those with a predicted score greater than the null relation) between each couple.

**Entity Resolution:** The previously extracted information constitutes a support for the entity resolution step. Each entity in the text is associated with an entity in the database, if possible. This allows to add new knowledge by completing the profile of the known entities or by creating new ones. In this pipeline, the entity resolution is solved by performing a search by mention and a selection by popularity. To each entity in the text, a list of KB entities is associated, that share the same type and at least one mention. In case the mentions do not return any results, an extended search is performed with the acronyms of these mentions. If no element is returned, the textual entity is added to the database. If several entities of the database match mentions of the textual cluster, a selection by popularity is applied, similarly to (Al-Badrashiny et al., 2017). The entity with the most occurrences, considering all mentions, is selected.

## 3.3 Coherence checks

The process described in the figure 2 concerns maritime events.

When a new event occurs, entities involved in maritime events are extracted from either the text or the AIS message: equipment, locations and organizations in our case. The KB is browsed, and a search is launched to check three conditions: whether there is an event of the same nature[2], involving the same ship(s) and occurring at the same date.

(i) If these three conditions are met, a similarity score $Sim_{wd}(E_i, E_j)$ is computed, using weights $w_x$ for each attribute $x$. These weights take into account the relations and attributes that are likely to embed misinformation or false information: if

---

[2]Ten natures of events are identified: seizure/arrest, collision, damage, sink, attack, aground, entrance (a harbor),leave (a harbor), transshipment and traffic.
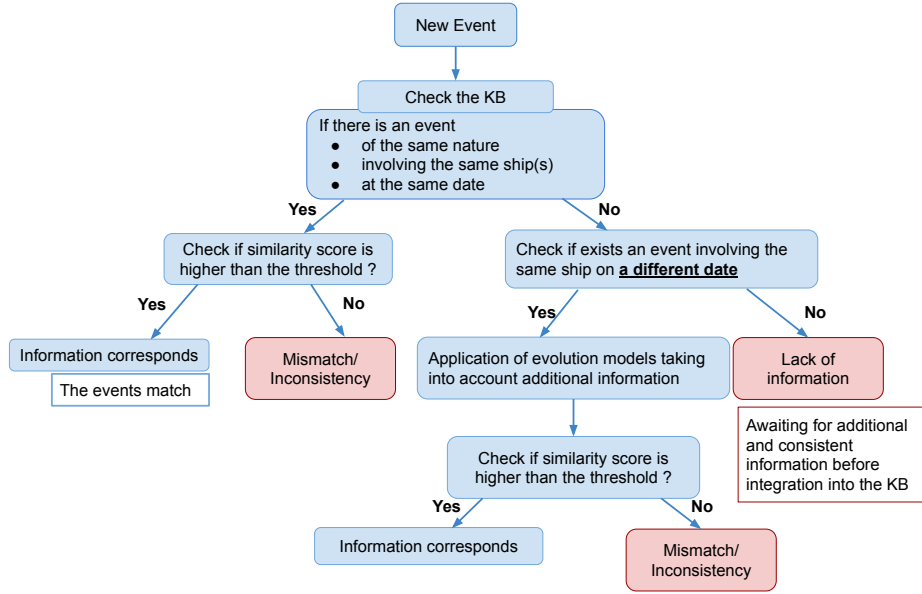
Figure 2: New event verification process

the location and/or the unit involved in the two events -to be compared- are different, there is a higher probability that one of these events contains misinformation. Attributes in the KB are compared one by one using the Jaro-Winkler distance (Jaro, 1989; Winkler, 1990). This score ranges between -1 and 1.

$$Sim_{wd}(E_i, E_j) = \sum_{\forall x \in S} w_x . d_{jw}(E_i, E_j) \quad (1)$$

where

- S ∈ [Org, Loc, Equipment, Pers]

- $w_x$ corresponds to the weight attributed to each entity.

- $d_{jw}(E_i, E_j)$ corresponds to Jaro-Winkler distance between event $E_i$ and event $E_j$.

A threshold is fixed to 0.25: if the similarity score is less than 0.25, the event is considered to hold disinformation. If the similarity score is higher than the threshold, that means that both events are considered coherent and the information brought by the new event corresponds to the information already in the KB.

(ii) If the three conditions are not satisfied, and an event involving the same ship on a *different date* already exists in the KB, additional information is considered by applying evolution models. Once evolution models are applied, the similarity score is computed again, and according to it, either the

two events match or an incoherence between the two events is spotted.

In the case where there are no events involving the same ship on a *different date* in the KB, there is a lack of information and the system cannot decide about coherence until more information is received.

**Evolution models**

Matching time-distant facts requires to consider a wider spectrum of evolution (e.g. for position) during the time difference. A simple similarity between the two facts would not be effective. Three evolving models for real-world application on ships have been identified and patented (Vasnier et al., 2022).

Each attribute in the KB (such as the name of the ship, its speed, location, etc.) is related to one of three types of evolution models, depending on the nature of attributes:

- constant model: constant attributes such as IMO (International Maritime Organisation) which is a unique identification number for ships are related to this type of model,

- predictable model: attributes that evolve over time such as the position, the direction or the speed of a ship are related to this type of model. The evolution of this kind of attributes is predictable with mathematical tools. Knowing the position of a ship and its heading direction, the geographical area in which the ship will be in a near future is predictable.

513

- **circumstantial model:** this type of model is the most complex to represent and to predict. It is related to events having attributes or relations which may change on rare occasions. The attributes related to the circumstantial model are subject to change with a specific and unpredictable event. In a maritime use case, the event could be the change of the captain, or further the purchase of a ship by another company.

The similarity score is computed as follows: $\forall p \in (E_i \cup E_j)$,

$$Sim_{evol}(E_i, E_j) = \frac{\sum (dist(p_{E_i}, p_{E_j}).\gamma_p)}{\sum (\gamma_p)} \quad (2)$$

where $\gamma_p$ denotes the confidence weight for each property $p$ of an event. $\gamma_p$ is the product of (a) the reliability of the sensor that collects information on $p$ and (b) the evolution uncertainty model of $p$. $\gamma_p$ is between 0 and 1. A weight of 1 is considered as a very reliable property and a weight of 0 means that we cannot trust this very uncertain property.

## 4 Performance analysis

### 4.1 K-pop Pipeline performance

**Setup:** To evaluate the information extraction pipeline, we focused on the proportion of information correctly extracted and aggregated from texts into a KB. In further detail, we computed a similarity score between a base populated by the evaluated system and the ground truth KB that we should obtain from a finite set of texts. For this purpose, we tested two scenarios, a *Warm-start* scenario which consists in populating an existing base and a *Cold-start* scenario in which we build a KB from scratch. To this end, we used the DWIE (Zaporojets et al., 2021) dataset. This dataset consists of 800 press articles in English, written and published by Deutsche-Welle. The textual level annotations of entities, their relations, their types and a unique identifier per entity at the inter-textual level allowed us to evaluate and compare our pipeline with the model proposed by (Zaporojets et al., 2021). The pipeline has been adapted to the ontology associated with the dataset and trained on the first 700 texts that constitute the train set. To measure a similarity score, we first align entities between the two KBs using the proportion of elements in common. The Hungarian algorithm (Kuhn, 1955) is then used to optimize this alignment, thus maximizing the average F1-score. Since the model introduced by

DWIE does not solve the entity resolution task, we use the same solution as the one in our pipeline.

**Results** The results in the table 1 illustrate the better performances compared with the DWIE model. Our K-POP pipeline shows up to a 2% improvement over the DWIE model in the *Warm-start* scenario. This shows that additional information extracted by our system contributes to a better linking with the existing content. The difference in results between the two types of scenarios shows the difficulty of populating an KB. However, our model shows a better resilience due to the linking by context approach.

| Model | F1 Cold-start | Warm-start |
|---|---|---|
| KBP | **76.1** | **72.1** |
| DWIE | 75.6 | 69.9 |

Table 1: F1 scores on the DWIE dataset.

Although there is still room for improvement, even more so in the case of the *Warm-start* scenario which shows the difficulty of populating an existing base, our IE (Information Extraction) solution can be considered for semi-automatic population.

### 4.2 Event consistency check

The aim is to ensure that the information extracted from the new event are coherent with those in the KB. If there is a contradiction with stored information in the KB, then an alert is raised. Figure 3 shows two events about the *Stena Impero* seizure in 2019. These events are extracted from two dif-
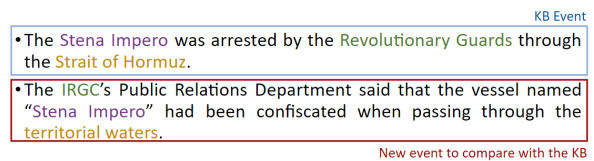


Figure 3: Example of two incoherent events.

ferent newspapers and they contain non-matching (incoherent) information. They are represented in the intelligent Knowledge Base as in figure 4.

The two events are compared based on the similarity score described in section 3.3.

Weights and threshold used for the determination of the similarity score are application-dependent; they rely on the sensibility of the end-user, since the results may vary depending on optimistic or
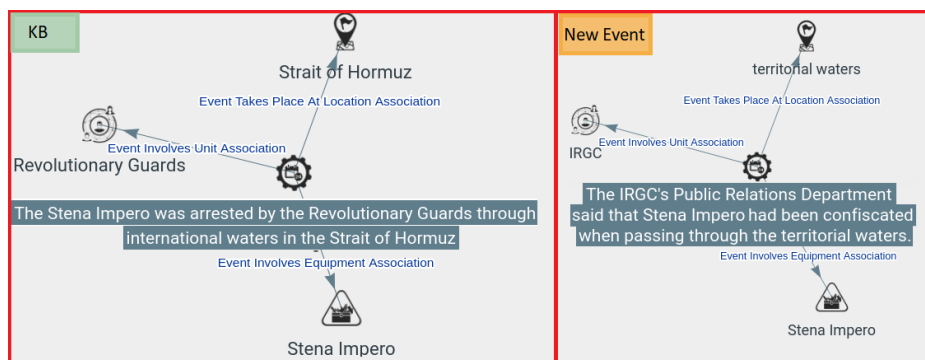
Figure 4: Representation of relation graphs of the event in the KB and the new event to evaluate

pessimistic assessment choices. They are currently defined based on end users application needs but can be further extended to allow for automated learning of these parameters. Note that the constitution of a training dataset for such a specific, unbalanced, high-risk problem is highly non-trivial. In this case, the new event holds disinformation since similarity score is -0.6 between these two events.

### 4.3 Evolution models in the real world

As an example, a newspaper may relate the following event, which UMBAR will need to compare with the existing events in the KB: *"On Saturday Stena Impero tanker had collided with a fishing boat, the Konarak, on its route."*

Extracted information from this event are in the table 2. We notice that there is a date -19 July

| Event | Stena Impero tanker had collided with a fishing boat, the Konarak, on its route. |
|---|---|
| Equipment | Stena Impero Konarak |
| Unit | - |
| Loc. | Bandar Abbas |
| Nature of event: | Collision |
| Date: | 19 july 2019 |

Table 2: Event and extracted information

2019- deduced from the publication date of the article containing the event.

This event conflicts with an event already present in the KB, indicating that *The konarak is moored in Turkey on the 13th of July*. Considering that this last event is correct, evolution models are used to perform the coherence check.
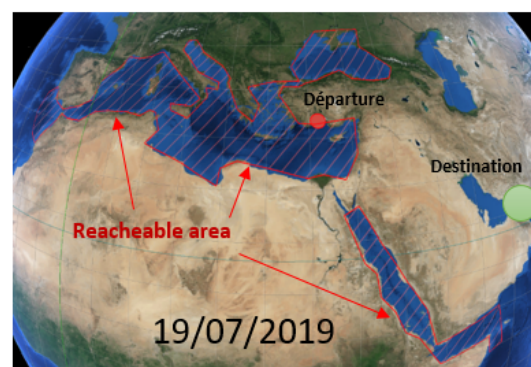


Figure 5: Reachable area from Turkey in five days

According to predictable models, a reachable area in five days (from 13 to 19 July) from Turkey is computed. The ship cannot be at the port of Bandar Abbas in Iran such a short time.

While these elements of evaluation are not provided as statistical measures, the specificity of the domain (high security risk along a low number of positive samples) makes it appropriate to evaluate the evolution models with an operational perspective.

## 5 Expert opinion and maritime security

Qualitative analyses on maritime security use cases are still on progress. For production-grade real word applications, the extracted information from AIS messages as presented in the figure 1 will be processed so as to check the coherence of information over time and to raise alerts in case suspect information is spotted.

Once combined, the aforementioned methods present features allowing to deal with actual data and information coming from heterogeneous sources, on a massive scale. It encompasses techniques to better correlate and assess information

with different timings, from both texts and sensor data domains, so that trusted KBs can be populated with a user-defined reliability.

Obviously, human intervention cannot be eliminated, but adding this process to ensure maritime security will help a lot, and will increase the performance of detecting false information or attempts to manipulate information.

# 6 Conclusion

This paper allowed us to present our semi-automatic end-to-end processing chain for information extraction and misinformation detection applied to a maritime surveillance use case. Maritime security being a central sector susceptible to false information leading to disastrous consequences. Although there is still room for improvement, our information extraction system and inconsistency detection provide support and alleviate the task of the operational staff in charge of monitoring. Future work will focus on improving the KBP pipeline to move towards fully automatic extraction, conducting an evaluation and further study on the detection of erroneous information.

The use of UMBAR in a representative setting is planned in order to evaluate the system and qualify it for its future operational deployment.

## Limitations

Building this system was nothing trivial. In our understanding the main challenges where to obtain data access, to chain very specialised artificial intelligence models, and to handle the iterations between the (machine) knowledge model, the customer expertise and the algorithms. We detail each of these challenges herein.

### Access to annotated data

Piracy and AIS spoofing are still too frequent, even though not frequent enough so as to result in the availability of datasets to train and evaluate an automatic system. The proposed approach mainly relies on subtask evaluation (notably on the information extraction steps). The Coherence Check is fully parameterizable in order to choose a sensibility to all possible variations. A stream of work concerning the automatic/statistic evaluation of the full pipeline is still going-on.

### Hyper-specialized AIs

Most of the substasks here are instantiated by trained modules, which inherently contain an adherence to the ontology used for labelling the training dataset. Information extraction from texts were fine-tuned for short pieces of news, and limited to English. This cuts off numerous relevant sources of information, typically from local newspapers anywhere on Earth.

### Handling business, ontology and algorithms together

The trend to fully automatize screening processes seems intuitive for many data scientists, but is actually not desirable for a security point of view: first, because the targeted elements are "black swans" which occur far too little in the training datasets, and more often than not, do not appear twice. Moreover, having too much confidence in the machine is clearly identified as a security risk, among other AI-system biases(Rastogi et al., 2020). Instead, the desired system should help the operator to handle more data about more incoming ships, and enabling them to focus on what is determining.

### Ethics Statement

Developing AI for security purposes always come with its ethical considerations. In this case, the system performs law enforcement and fight against piracy, which are commonly assessed as noble, ethical deeds. As the application specifically targets maritime trafficking, the risk of misuse is reduced (i.e. it cannot be used to target individuals).

This system relies on third party sources of data. As a consequence, the data processing roles are clearly and contractually established between the providers and the customer of this system, decreasing privacy risks. No personal data is required by the system; personal *public* data may be handled from the press and from the AIS information (typically, the name of the captain).

The final result of the system is to gather and aggregate a complete picture of the risk level of a ship, to help an operator. The system may be used to prioritize the effort to review the documents and cargo of a ship, but cannot be used to authorize or forbid a ship's entry – this remains the decision of the operator.

# References

Mohamed Al-Badrashiny, Jason Bolton, Arun Tejasvi Chaganty, Kevin Clark, Craig Harman, Lifu Huang, Matthew Lamm, Jinhao Lei, Di Lu, Xiaoman Pan, et al. 2017. Tinkerbell: Cross-lingual cold-start knowledge base construction. In *TAC*.

Houssein Alaeddine and Cyril Ray. 2022. A hybrid artificial intelligence system for securing a maritime zone based on historical and real-time data analysis. In *OCEANS 2022, Hampton Roads*, pages 1–8. IEEE.

Marco Balduzzi, Alessandro Pasta, and Kyle Wilhoit. 2014. A security evaluation of ais automated identification system. pages 436–445.

Rajarshi Bhowmik and Gerard de Melo. 2020. A joint framework for inductive representation learning and explainable reasoning in knowledge graphs. *CoRR*, abs/2005.00637.

Ziyang Chen, Xiang Zhao, Jinzhi Liao, Xinyi Li, and Evangelos Kanoulas. 2022. Temporal knowledge graph question answering via subgraph reasoning. *Knowledge-Based Systems*, 251:109134.

Claudia d'Amato, Pierpaolo Masella, and Nicola Fanizzi. 2022. An approach based on semantic similarity to explaining link predictions on knowledge graphs. In *IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, WI-IAT '21, page 170–177, New York, NY, USA. Association for Computing Machinery.

Vladimir Dobrovolskii. 2021. Word-level coreference resolution. *arXiv preprint arXiv:2109.04127*.

Jacques Everwyn, Bruno Zanuttini, Abdel-Illah Mouaddib, Sylvain Gatepaille, and Stephan Brunessaux. 2019. Achieving maritime situational awareness using knowledge graphs: a study. In *1st Maritime Situational Awareness Workshop (MSAW 2019)*.

Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206.

Matthew A. Jaro. 1989. Advances in record-linkage methodology as applied to matching the 1985 census of tampa, florida. *Journal of the American Statistical Association*, 84(406):414–420.

Harold W Kuhn. 1955. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Maxime Prieur., Cédric Mouza., Guillaume Gadek., and Bruno Grilheres. 2023. Evaluating and improving end-to-end systems for knowledge base population. In *Proceedings of the 15th International Conference on Agents and Artificial Intelligence - Volume 3: ICAART,*, pages 641–649. INSTICC, SciTePress.

Charvi Rastogi, Yunfeng Zhang, Dennis Wei, Kush R Varshney, Amit Dhurandhar, and Richard Tomsett. 2020. Deciding fast and slow: The role of cognitive biases in ai-assisted decision-making. *arXiv preprint arXiv:2010.07938*.

Noureddine Seddari, Abdelouahid Derhab, Mohamed Belaoued, Waleed Halboob, Jalal Al-Muhtadi, and Abdelghani Bouras. 2022. A hybrid linguistic and knowledge-based analysis approach for fake news detection on social media. *IEEE Access*, 10:62097–62109.

Kilian Vasnier, Sylvain Gatepaille, and Valerian Justine. 2022. Method and system for merging information. Patent No. US20220374464A1. https://patents.google.com/patent/US20220374464A1.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

William Winkler. 1990. String comparator metrics and enhanced decision rules in the fellegi-sunter model of record linkage. *Proceedings of the Section on Survey Research Methods*.

Klim Zaporojets, Johannes Deleu, Chris Develder, and Thomas Demeester. 2021. Dwie: An entity-centric dataset for multi-task document-level information extraction. *Information Processing & Management*, 58(4):102563.

Jiasheng Zhang, Shuang Liang, Yongpan Sheng, and Jie Shao. 2022. Temporal knowledge graph representation learning with local and global evolutions. *Knowledge-Based Systems*, 251:109234.

Wen Zhang, Bibek Paudel, Wei Zhang, Abraham Bernstein, and Huajun Chen. 2019. Interaction embeddings for prediction and explanation in knowledge graphs. *CoRR*, abs/1903.04750.

Wenxuan Zhou, Kevin Huang, Tengyu Ma, and Jing Huang. 2021. Document-level relation extraction with adaptive thresholding and localized context pooling. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 14612–14620.