MILITARY UNIVERSITY OF TECHNOLOGY FACULTY OF CYBERNETICS INSTITUTE OF INFORMATION SYSTEMS



Identification of hidden semantic relations in texts using relations' patterns.

Dissertation

submitted by

mgr inż. Michał Gałusza

supervised by: prof. Andrzej Walczak Warsaw, 2024

Institute of Information Systems Military University of Technology

Declaration

.....

.....

Warsaw, 31st Jan 2024

Your signature

Acknowledgements

Completing this PhD thesis has been a challenging and rewarding journey, and it would not have been possible without the support and guidance of many individuals. I am deeply grateful to all those who have contributed to this work and supported me throughout this process.

First and foremost, I would like to express my heartfelt gratitude to my supervisor, Prof. Andrzej Walczak, for his guidance, encouragement, and trust in my exploration of knowledge avenues in risk analysis and natural language processing areas. Although most of them often led to dead ends, which required me to backtrack and redo my work numerous times, Your patience has been instrumental in shaping this research and helping me navigate the complexities of the academic world.

I extend my thanks to my colleagues and friends in the Cybernetics Faculty, Institute of Information Systems, for their camaraderie and support. I especially thank Prof. Andrzej Chojnacki for the stimulating discussions during the seminar sessions.

I want to acknowledge the administrative and technical staff at the Military University of Technology for their assistance and support. Your help in navigating the various bureaucratic and logistical hurdles has been invaluable.

I am indebted to my family for their unwavering support and encouragement, to my wife, Marta, and my kids, Maks and Margaret, your patience, understanding, and encouragement have been my anchor throughout this journey. To my parents, Danuta and Leszek, thank you for instilling in me the value of education and your unconditional love and support. I will now have more time for my children, .

Thank you all for your support and belief in me. This thesis is a testament to the collective effort and encouragement I have received from each of you.

Podziękowania

Ukończenie tej pracy doktorskiej było wymagającą i satysfakcjonującą podróżą, która nie byłaby możliwa bez wsparcia i wskazówek wielu osób. Jestem głęboko wdzięczny wszystkim, którzy przyczynili się do tego dzieła i wspierali mnie w trakcie tego procesu.

Przede wszystkim chciałbym wyrazić moje serdeczne podziękowania mojemu promotorowi, prof. Andrzejowi Walczakowi, za jego przewodnictwo, zachętę i zaufanie do mojej eksploracji obszarów analizy ryzyka i przetwarzania języka naturalnego. Chciaż większość z nich często prowadziła do ślepych zaułków, i wiele razy wymagało to ode mnie powrotu i ponownego wykonania pracy to pańska cierpliwość była kluczowa w kształtowaniu tych badań i pomocy w nawigowaniu po zawiłościach świata akademickiego.

Serdecznie dziękuję moim kolegom i przyjaciołom z Wydziału Cybernetyki, Instytutu Systemów Informacyjnych za towarzystwo i wsparcie. Szczególne podziękowania dla prof. Andrzeja Chojnackiego za inspirujące dyskusje podczas sesji seminaryjnych.

Chciałbym podziękować administracji i pracownikom technicznym na Wojskowej Akademii Technicznej za ich pomoc i wsparcie. Wasza pomoc w poruszaniu się po różnych biurokratycznych i logistycznych przeszkodach była nieoceniona.

Jestem wdzięczny mojej rodzinie za ich niezmienne wsparcie i zachętę. Mojej żonie, Marcie, dzieciom Maksymilianowi i Małgorzacie, dziękuję za cierpliwość, zrozumienie i wsparcie, które były moją kotwicą przez całą tę podróż. Moim rodzicom, Danucie i Leszkowi, dziękuję za wpajanie mi wartości edukacji oraz za bezwarunkową miłość i wsparcie.

Dziękuję wszystkim za wsparcie i wiarę we mnie. Ta praca jest świadectwem zbiorowego wysiłku i zachęty, które otrzymałem od każdego z Was.

Abstract

The thesis proposes a solution to the problem of recognizing and modeling risk transmission contained in documents describing a system's operation or failure. Relationships are determined in the context of entire documents and go beyond the currently dominant approach of identifying relationships within a single sentence and using classifiers trained on collected dedicated training examples.

The problem being addressed is significant because information about the flow of threats can create complex interactions between system elements described in such a way that they may be scattered throughout the document and even between sources. There is generally a lack of dedicated training sets for classifiers, and existing solutions are limited to selected areas, such as railways, or description formats, such as HAZOP or FMEA.

The proposed solution involves a gradual decomposition of descriptions. The examined text is decomposed into a Semantic Frames Graph (SFG) in the first step. In the second step, the pattern of threat propagation relationships is used to recognize propagation. Recognized propagations are stored in an Intermediate Relationship Graph (IRG). In the final step, propagations are aggregated into the form of an Asset-Vulnerability-Hazard (A-V-H) graph, which allows for a network analysis of the risk contained in the description of the operation of a given system.

The proposed approach allows for modeling risk propagation without needing a dedicated relationship detection mechanism, as this method is based on verbalizing the relationship pattern. Another reason for eliminating dedicated classification is the extension of pattern analysis to analyze the dialog coherence in the path between nodes in the SFG graph. The detection results obtained by combining both methods are verified using current language models (Large Language Models such as chatGPT) and prompt engineering. The threshold above which relationships are accepted is a solution to the multi-criteria optimization task.

Overall, this work presents a new method for detecting relationships and its application in risk analysis. It also explores the potential of semantic pattern methods, dialogic coherence, and prompt engineering in constructing a network risk model, which facilitates modeling complex threat propagation dependencies.

Streszczenie

Rozprawa proponuje rozwiązanie dla problemu rozpoznawania i modelowania transmisji ryzyka zawartego w dokumentach opisujących działanie systemu lub opisujących jego awarię. Relacje wyznaczane są w konteście całych dokumentów i wykraczają poza aktualnie dominujące podejście wyznaczania relacji w ramach jednego zdania oraz przy użyciu klasyfikatorów wytrenowanych na zebranych dedykowanych przykładach trenujących.

Rozwiązywany problem jest istotny jako, że informacje o przeływie zagrożenia mogą tworzyć skomplikowane interakcje pomiędzy elementami systemu opisanymi w taki sposób, że mogą być rozproszone po całym dokumencie a nawet pomiędzy źródłami i na ogół brakuje dedykowanych zbiorów trenujących klasyfikatory a istniejące rozwiązania są ograniczone do wybranych obszarów np.: kolei, lub formatów opisow np.: HAZOP lub FMEA.

Proponowane rozwiązanie zakłada stopniową dekopmozycję opisów. W pierwszym kroku badany tekst dekomponowany jest do postaci Grafu Ramek Semantycznych (ang. Semarntic Frames Graph, SFG). W drugim kroku, wzorzec relacji propagacji zagrożenia używany jest do rozpoznania propagacji. Rozpoznane propagacje zapisywane są w grafie Pośrednich Relacji Semantycznych (ang. Intermediate Relathionship Graph, IRG). W ostatnim kroku, propagacje są agregowane do postaci grafu Zasób-Podatność-Zagrożenie (ZPZ) (ang. Asset-Vulnerability-Hazard, A-V-H), który pozwala na sieciową analizę ryzyka zawartego w opisie dzialania danego systemu.

Zaproponowane podejście pozwala na modelowanie propagacji ryzyka bez konieczności stosowanie dedykowanego mechanizmu detekcji relacji jako, że metoda ta opiera sie na werbalizacji wzorca relacji (ang. pattern verbalization). Drugim powodem, który pozwana na eliminację dedykowanej klasyfikacji jest rozszerzenie analizy wzorca o analizę spójnościa dialogowej w scieżce pomiedzy węzłami w grafie SFG. Wyniki detekcji uzyskanych poprzez zestawienie obu metod weryfikowane są poprzez wykorzystanie aktuanych językowch modeli generatywnych (Large Language Models np.: chatGPT) oraz inżynierii podpowiedzi (ang. prompt engineering). Próg powyżej którego relacje są akceptowane jest rozwiązaniem zadania optymalizacji wielokryterialnej.

Ogólnie nieniejsza praca przedstawia nową metodą detekcji relacji i jej zastosowanie w obszarze analizy ryzyka. Przedstawia potencjał metody wzorców semantycznych, spójności dialogowej oraz inżynierii podpowiedzi w konstruowaniu sieciowego modelu ryzyka, który ułatwia modelowanie złożonych zależności propagacji zagrożenia.

Acronyms

- **RA** Risk Analysis
- LM Language Model
- **LLM** Large Language Model
- **RM** Risk Management
- SFG Semantic Frames Graph
- SRL Semantic Role Labelling
- IRG Intermediate Relationships Graph
- A-V-H Asset-Vunerability-Hazard
- **OIE** Open Information Extraction
- KG Knowledge Graph
- **KA** Knowledge Acquisition
- KAP Knowledge Acquisition Pipeline
- **SN** Semantic Network
- PHA Pre-Hazard Analysis
- HAZOP Hazard and Operability Study
- **OIE** Open Information Extraction
- TL Tranfer Learning
- **NER** Named Entity Recognition
- NLI Natural Language Inference
- **RE** Relationship Extraction

Contents

Ac	cknow	wledge	ments	i	
Po	odzię	kowan	ia	ii	
Al	ostra	ct		iii	
St	reszc	zenie		iv	
Ac	crony	ms		v	
1	Intr	oducti	on	1	
	1.1	Thesi	s structure	5	
	1.2	Motiv	rations	6	
		1.2.1	Risk Analysis	6	
		1.2.2	Risk - Asset - Vulnerability Dilemma	8	
		1.2.3	Large Language Models	9	
		1.2.4	Validation and Explainablility	12	
		1.2.5	Goal of Reasearch	13	
2	Rela	ated W	ork	15	
	2.1	Ontol	ogy in Risk Analysis	17	
		2.1.1	Direct Application of Reasoning on Ontologies in Risk Analysis	18	
		2.1.2	Application of Ontology in the Knowledge Graph Construction	23	
	2.2	2 Knowledge Graphs and Network Representation of Risk			
		2.2.1	Representing HAZOP Safety Reports as Knowledge Graphs	28	
		2.2.2	Representing Free Text Risk-Related Narratives as Knowledge Graphs .	29	
	2.3	Sumn	nary	30	
3	Rele	Relevant Natural Language Processing Techniques			
	3.1	Entity	Recognition	33	
	3.2	Relati	onship Extraction	35	
		3.2.1	Intra-Sentence Relationship Extraction	35	
		3.2.2	Inter-Sentence Relationship Extraction	37	
		3.2.3	Transformer-based Relationship Extraction	43	

	3.3	Textual Entailment 46
	3.4	Semantic Frames And Semantic Role Labeling
	3.5	Summary
4	Proj	posed Solution 53
	4.1	Asset-Vunerability-Threat Triplet and A-V-H Graph
	4.2	Problem Statement 54
	4.3	Proposed Solution Architecture
	4.4	Sentence Decomposition and Semantic Role Labelling
	4.5	Semantic Frames Graph
	4.6	Relationship Extraction and Semantic Pattern
		4.6.1 Single Frame Relation Extraction
		4.6.2 Two-frame Relation Extraction
		4.6.3 Modified Dialog Coherence Function
		4.6.4 Combined Entailment Function
		4.6.5 Multiple Templates
	4.7	Intermediate Relationship Graph
	4.8	Asset-Vulnerability-Hazard Graph
	4.9	Validation and Explainability
		4.9.1 Validation and Threshold Calculation
		4.9.2 Explainablity
	4.10	Summary
5	Res	ults 78
	5.1	Knowledge Acquisition Pipeline
	5.2	Intra-Sentence Relation Detection
	5.3	EMARS Report Example
	5.4	Impact of Large Langauge Models
	5.5	Additional Examples
		5.5.1 Financial Scenario
		5.5.2 Medical Scenario
	5.6	Summary

6 Conclusions

2

Biblio	graphy	7	113
.1	Com	putation Complexities Calculatation	124
	.1.1	Naive Relation Validation	124
Appen	dix		124

Chapter 1

Introduction

Cambridge English Dictionary defines the adjective "hidden" as «not easy to find». Although "not easy" is connected with "difficult", in the domain of risk analysis, it should instead refer to the fact of not being directly detectable. The «hidden» nature of risk relations results from human language's flexibility in expressing the descriptions of risk. In addition, the semantics' of risk is characterized by extreme contextuality.

The «hidden» aspect of risk interactions within a system often presents challenges in modeling and analysis. At the heart of this complexity lies the intricate interplay between various factors, each exerting its influence and contributing to the overall risk propagation. These interactions are seldom linear or predictable, making them difficult to capture within traditional risk models. Moreover, the interconnected nature of modern systems introduces a web of dependencies and feedback loops, amplifying the potential for cascading effects and undetected consequences. Despite these challenges, understanding and effectively managing complex risk interactions are essential for safeguarding against systemic failures and disruptions. This requires a holistic approach that integrates diverse perspectives, leverages advanced modeling techniques, and embraces uncertainty to effectively uncover hidden vulnerabilities and mitigate emerging threats. Through such combined approaches, we can navigate the intricate web of risk interactions and effectively safeguard against hazards.

Linguistic aspects of relationship extraction focus on identifying and categorizing semantic connections between entities within and across sentences, leveraging both syntactic and semantic analyses. Intra-sentence relations involve identifying relationships within a single sentence, which relies heavily on syntactic parsing to understand grammatical structures and dependencies, such as subject-object relationships, as in "The scientist discovered a new element." Semantic roles further clarify these relationships by specifying the function of each entity, like identifying "the scientist" as the agent and "a new element" as the object. Such an approach simplifies the detection as the scope of analysis is limited to the sentence. Inter-sentence relations extend this task to link entities across multiple sentences, requiring discourse analysis to resolve co-references and maintain entity consistency. For example, in the text "Marie Curie conducted pioneering research on radioactivity. She was the first woman to win a Nobel Prize," co-reference resolution identifies "She" as referring to "Marie Curie," thus connecting the two sentences and establishing the relation between "Marie Curie" and "Nobel Prize". This comprehensive approach, integrating intra- and inter-sentence analyses, ensures a robust understanding of entity relationships, essential for applications in information retrieval and knowledge base construction.

Formulating the relationship extraction task involves defining it in a structured manner, incorporating the identification of entities and the semantic relations between them within a given text.

1. Input:

- A corpus of text *T*, which is a collection of sentences $S = \{s_1, s_2, ..., s_n\}$.
- Each sentence s_i consists of a sequence of words $\{w_1, w_2, \ldots, w_m\}$.

2. Named Entity Recognition (NER):

Identify and classify named entities within the text. The output is a set of entities $E = \{e_1, e_2, ..., e_k\}$, where each entity *e* is associated with a type (e.g., Person, Organization, Location).

3. Candidate Relation Identification

Generate pairs of entities (e_i, e_j) from the set *E*. Each pair is considered a candidate for potential relationships.

4. Feature Extraction:

Extract linguistic features from the text to help determine the relationship between entities. These features can include:

- Syntactic features: Dependency paths, part-of-speech tags, parse trees.
- Semantic features: Word embeddings, semantic roles, named entity types.
- Contextual features: Surrounding words, co-reference chains, sentence position

5. Relation Classification

For each candidate pair (e_i, e_j) , classify the semantic relation r between the entities. This is often treated as a multi-class classification problem where the output is a predefined set of relation types $R = \{r_1, r_2, ..., r_l\}$, including a "no relation" class if applicable.

6. Output

A set of triples $\{(e_i, r, e_j)\}$ where each triple represents a relationship *r* between entities

 e_i and e_j .

By structuring the relationship extraction task, we can systematically identify and classify relationships between entities within and across sentences, leveraging various linguistic features and techniques. The approach is successful in scenarios where the relations are expressed **explicitly**. Consider the example: s = "Marie Curie discovered radium in 1898.". In this case, the approach would address the detection of "radium discovery" is the following manner:

1. Named Entities

- "Marie Curie" (Person) - "radium" (Substance)

2. Candidate Relation

("Marie Curie", "radium")

3. Feature Extraction

Syntactic dependency: "Marie Curie" (subject) -> "discovered" (verb) -> "radium"
 (object) - Semantic roles: "Marie Curie" (Agent), "radium" (Object)

4. Relation Classification

Classify the relation between "Marie Curie" and "radium" as "discovered"

5. **Output** The relationship triple: ("Marie Curie", "discovered", "radium").

The nature of risk-related relations challenges the standard approach to RE in two ways. First, these relations are **implicit**, meaning there is no direct mention of specific risk-related relations in the text. Consider the sentence: "ATF is a type of aviation fuel designed for use in aircraft-powered gas-turbine engines. If these supercooled droplets collide with a surface, they can freeze and may result in blocked fuel inlet pipes". An average human reader easily spots the impact of a droplet on the engine even though there is no **direct** mention of it. This case contrasts with Marie Curie's radium discovery in the previous example.

Second, the named entities should be risk-specific, which means that instead of general categories "Person," "Organization," or "Location," we are more interested in related categories such as "Asset," "Hazard," or "Barrier" Consider a sample sentence: "A hunter shot a raging bear". Both "a hunter" and "a bear" are *Threats* depending on the context. For an object "bear, "hunter" is a *Threat* or *Hazard* as he shot it eventually; however, for an object 'hunter," "bear" is a *Threat* therefore, it was shot.

Contextuality is even more visible if we expand the scope. Let's consider another sentence: "A hunter shot a raging bear attacking a woman". Within this single sentence, object "hunter" should be assigned two roles simultaneously: *Threat* from the bear's perspective and *Savior* or *Barrier* from the woman's. Therefore, which concept should represent the word "hunter" in this sentence?

In the risk domain, such contextual situations are not uncommon and require a contextual approach to detect them correctly. For example, the contextual role of the package (Fig. 1.1) depends on whether we consider a human underneath - in this case, it will belong to class *Hazard*. However, given that the package is valuable, it will be considered as *Asset*, which *Vaulnerability* would be a line carrying it and *Hazard* an event of the line snapping.

The contextuality challenges a classical knowledge graph construction as it relies on the notion of *Concept*. The hierarchy of concepts and the relations between them is the foundation of the representation - the model of the domain of interest. That organized structure forms an *Ontology* a foundation for *Inference*. In the classical approach, a concept is a component of human thought and is the thinking unit that refers to objective things and their peculiar properties. A concept's formation is a procedure with the direction "from special to general". Considering various objects that are "special" cases, one determines a "general" set of properties that form the concept. This implies that we can define the concepts only through their properties and how linguistic expressions of concepts exist within the narrative. As with the word "apple," we can associate the information related to its shape, color, taste, and the context in which it usually appears in any narrative. We can observe that other words, such as "peach" and "banana," share the same linguistic properties; therefore, they are similar. All of these allow classifying "apple", "peach", and "banana" to the concept of "fruit".

Using a similar approach to identifying, for example, the *Hazard* concept would require enumerating its features, effects, and linguistic expressions to detect them in the parsed descriptions. Due to the unmanageable number of combinations, such an approach is impossible as the concept of hazard manifests itself in various domains differently. For example, in the medical domain, the hazard manifests through adverse drug effects, impact on the organs, or general deterioration of patients' medical conditions. In the financial realm, risk will be manifested through capital loss. In software engineering, through a data breach or unexpected system malfunction.

The thesis aims to extract risk-related interactions from the text. It uses a specific risk structure, a triple Asset-Vulnerability-Threst [1], that organizes and constrains how the interactions are identified. The approach proposed relaxes the problem of direct detection of risk-related concepts and formulates the methodology, which allows constructing the



Figure 1.1: Risk Contextuality

comprehensive representation of risk interactions in the form of a specific Knowledge Graph called Asset-Vulnerability-Hazard graph (pol. graf Zasób-Podatość-Zagrożenie) [1].

1.1 Thesis structure

The thesis is structured as follows:

• Chapter 1 – Introduction

This chapter explains the motivations behind the research. It also outlines challenges in Risk Analysis and requirements for a solution to construct a comprehensive risk representation.

• Chapter 2 – Related Work

The chapter summarizes related work in network representation of risk and discusses the limitations of current methods from the representation and text processing perspective.

Chapter 2 – Relevant Natural Language Processing Techniques

The chapter discusses relevant Natural Language Processing techniques used to construct the network representation of risk. It identifies the bottlenecks in the general Knowledge Acquisition Pipeline, which is used to parse and transform the textual representation of risk, especially in entity recognition and relationship extraction areas.

• Chapter 4 – Proposed Solution

This chapter discusses a proposed solution for a risk modeling system. It explains

the main challenge in the naive approach to solving required triple identification based on entity recognition. It provides a solution via analysis of risk propagation and describes the main solution concepts: Semantic Frame Graphs, Dialog Coherence, and Intermediate Relationship Graphs. It formulates a multi-objective optimization to establish the threshold on dialog coherence scores on verbalizing the risk-relevant relationship templates.

• Chapter 5 - Results

This chapter presents the results, provides insights into how current LLM is used in validating relations, and explains how the solution can be used in domains other than risk.

Chapter 6 - Conclusions

This chapter concludes the dissertation and outlines future work.

1.2 Motivations

1.2.1 Risk Analysis

Critical factors support the continuous improvement of Risk Analysis and Risk Management methods. First, as we are faced with a "fast pace of technological change," "new types of hazards," and "increasing complexity and coupling." [2], it is necessary due to the growing complexity of new systems, objects, and processes in which risk-related interactions are increasingly difficult to find and model. Second, risk management must be an integral part of the overall management process, describing risk interactions in a meaningful and standardized way across system elements to allow for informed decisions on risk-preventing strategies. Third, legal regulations already impose risk management strategies on corporations, i.e., Seveso Directives.

However, few structured data sources collect risk interaction comprehensively, even though collecting, analyzing, and storing data relating to accidents and incidents, given the regulations, is mandatory in some industries. There are just a few examples of official government-managed semi-structured repositories that are textual and, therefore, do not allow systematic analysis without additional transformations:

• *Nuclear Power Industry.* In this industry, the data collection is rooted in the International Convention on Nuclear Safety. According to this convention, each contracting party commits to taking the appropriate steps to ensure that: «incidents significant to safety are reported in a timely manner by the holder of the relevant license to the regulatory body; [and that] programs to collect and analyze operating experience are established, the results obtained and the conclusions drawn are acted upon and that existing mechanisms are used to share important experience with international bodies and with other operating organizations and regulatory bodies» [3]

- *Aviation*. According to EU directive 2003/42/EC on "Occurrence reporting in civil aviation" data related to all civil aviation incidents and accidents must be collected, reported, and analyzed. The organization European Co-ordination Centre for Accident and Incident Reporting Systems (ECCAIRS) has been established to «assist national and European transport entities in collecting, sharing and analyzing their safety information to improve public transport safety.»
- *Process industries covered by Seveso II Directive.* Companies in Europe that comply with the Seveso II directive must collect and report data in a specified format to the national authorities and the eMARS database.

Understanding the connections between components, hazards, and consequences in the system's domain is a key element in reasoning about the propagation of hazards contained in textual risk repositories. Methodologies exist to construct representations focusing on specific risk modeling approaches, namely qualitative and quantitative risk assessment methods.

In reality, however, a limited model is usually constructed through an iterative, timeconsuming process involving subject matter experts with a focus on a selected area of operation of the system [4], [5]. This leads to a situation where most of the risk-related data is written down and stored in various descriptions, either of the failure events, i.e., railway accident report [6], or "near-misses" reports in industrial cases [7] or as descriptions of safety operations of the given system. Such representation restricts the possibility of analyzing the risk interactions comprehensively, as documents need to be read and interpreted by experts.

This defines the first set of requirements for building comprehensive risk representation. Such a system shall be able to consume information written in natural language, normalize it, and store it in a format allowing standardized analysis.

Another critical factor in risk prevention is the subjective nature of risk itself. It is influenced by a broad set of phenomena beyond the mere technical conception of risk as a combination of accident scenarios, probabilities, and adverse outcomes. Excessive reliance

on subject matter experts is a risk of the risk analysis methods.

The cognition bias of experts performing the risk analysis, regardless of methodology, given the system's complexity and time or budget constraints, is responsible for underestimating or even omitting scenarios that lead to a catastrophe. An excellent example of such accidents would be "Herald of Free Enterprise" or "Jan Heweliusz" - roll-on roll-off car and passenger ferries capsizing. In the first case, missing bow door indicators allowed the ship to depart with the bow doors unlocked. In the second, an inefficient weight-balancing mechanism was inadequate for weather conditions.

Therefore, the second set of requirements is to mitigate subjective risk perception. The system shall be constructed so that a clear template related to risk propagation should be used against the set of documents while searching for risk relations. The role of experts shall be switched from performing complete risk analysis manually to collecting relevant documentation and validating the results of risk detection. The system shall automatically ingest and transform the documents to normalized representation incrementally so the risk representation is augmented when new facts arrive.

The final requirement relates to representing risk interaction in the modeled system. The natural candidate would be the network model of risk interaction. Its advantage is the simplicity of interpretation, which means that the graphical form is understandable also for those not involved in construction. In the case of risk assessment, the undoubted advantage of the network model is the ability to visualize the links between the effects of threats, which is an initial step towards more complex quantitative risk estimation as Bayes Nets [1]. Defining dependencies directly in the Bayesian network is troublesome as it requires the decomposition of risk interaction and then estimating conditional risk probabilities [1]. Therefore, developing methods for building network security models may prove to be a foundation for cost-effective ways of security analysis.

1.2.2 Risk - Asset - Vulnerability Dilemma

The contextuality of risk means that the same element can belong to all classes. For example, "engine" is an asset impacted by the "droplet" risk in "fuel". However, the "airplane" is an *Asset* too that is impacted by the "engine" as the its flying capability relies on it. In this case, "engine" is the airplane's *Vulnerability*. Therefore, "engine" must be assigned two concepts simultaneously.

The current class detection, Named Entity Recognition (NER), in NLP pipelines relies on a text span classification approach, in which both span and the class are detected [8]. The approach is limited in two ways. First, it is limited by the selection of span sizes to capture the interaction, and second, by lack of referring between entities. In the second case, the A-V-T triple would require the classifier to assign multiple classes to the same span, denoting either entity depending on which element is considered an *Asset*. This means that variant entity classification is required in a single context, which is impossible in current NER solutions.

1.2.3 Large Language Models

Late rapid progress on Large Language Models (LLMs) was initiated with the publication of the Transformer architecture [9]. Two elements are behind LLMs' success. First is the attention mechanism, which allows weighted access to the fragments of the context. The other is the model's architecture, which enables easy expansion of the model's parameter space.

With increased computing and training resources, LLMs have demonstrated increased semantic capabilities, making perfect use of both the attention mechanism and its scalable architecture. Since Alan Turing's seminal paper on "Computing Machinery and Intelligence" and his famous Turing Test, we have progressed to the state of the art, which sparks ongoing discussions on threats posed by the uncontrolled growth of capabilities of such models [10]. Such discussions are academic no more, and on 31st October 2023, the UK Prime Minister, Mr. Rishi Sunak, hosted the first global summit on the risks associated with artificial intelligence.

The capabilities of the current general language models, called *foundation models* [11], are, in fact, staggering. However, they are limited in several aspects.

First, the performance is a function of their parameters [12]. A perfect example explaining the improvement with the increase of model parameters is a question answering where the zero-shot setting comes close to the current state-of-the-art performance of the fine-tuned models [13] (Fig. 1.2).

Increased parameter space comes with computing requirements. The scaling is visible if we compare models' training compute power utilization [13] (Fig. 1.3).

Lastly, significant parameter space requires a significant amount of data. OpenAI's GPT-3 model was trained on a filtered Common Crawl dataset, WebText, two internetbased books corpora, and English Wikipedia. All three elements, significant parameter space, compute resources, and abundant training data, enable identifying representations of the knowledge encoded in corpora at scale. Still, the model learns from the data it



Figure 1.2: Question Answering Performance [13]



Figure 1.3: Model Training Compute Cost Comparison [13]

used; Therefore, OpenAI had to undertake a deliberate training strategy to counter data contamination [13].

Due to data requirements, only a few areas have enough textual resources to train the dedicated LLMs. In the judicial domain, LLMs encode unstructured textual resources that comprise the legal system. The reason why specific training is required is that the nature of the language in the domain deviates from the 'common' language used daily. A very good example of such a case is the word "consideration", which in general English means "the act of thinking about something carefully" and in legal terms, it is "to describe the benefit each party to a contract receives". This is often payment in exchange for goods or services. In the judicial domain, LLMs can perform, among others, the following tasks [14]:

- they can quickly extract key points from legal documents, combine them with the judgment outcomes, and generate concise and accurate case summaries,
- they can generate a draft legal document that complies with legal standards,
- through interactive Q&A sessions with users, they can provide convenient and efficient legal consultation services, while also reducing the workload of professional lawyers,
- they can summarize and extract the key features of a given case, which can contain significant legal documentation.

In the financial area, LLMs are trained in specialized financial textual resources to perform [15]:

- Regulatory Compliance to assists financial institutions in analyzing and interpreting complex regulatory documents, compliance requirements, and legal agreements.
- Investment Research and Due Diligence to help analyze vast amounts of financial reports, company filings, analyst research, and news articles to identify investment opportunities and conduct due diligence on potential investments.

By connecting historical textual data and financial data, LLMs can perform rudiment financial analysis such as:

- Market Analysis and Forecasting in which LLMs analyze historical financial data, market trends, and economic indicators to generate insights and forecasts,
- Fraud Detection in which LLMs can assist in detecting fraudulent activities, suspicious transactions, and potential money laundering activities by analyzing textual data such as transaction records, customer communications, and public records.

Last but not least is the medical domain. It has a portfolio of dedicated language models, including MedPalm and MedPlam2 [16], which encode medical knowledge and have been

developed specifically for medical question-answering tasks. However, potential use cases for such models are more advanced and apply to areas such as:

- Clinical Trial Matching and Recruitment where LLMs can help match patients, based on their medical history, to relevant clinical trials,
- Healthcare Chatbots and Virtual Assistants where the chatbots and virtual assistants can be used to interact with patients to schedule appointments, answer medical questions, provide medication reminders, and offer health coaching,
- Integration of Electronic Health Records (EHR) where LLMs can improve the efficiency and accuracy of analysis of EHR by interpreting and extracting relevant information from unstructured clinical notes, physician narratives, and patient histories.

These models share common traits: they are significantly founded and have enough good-quality data for training. The risk domain has incomparably smaller data sets for several reasons:

- textual resources are scarce as once an adverse event occurs, corrective measures are taken to prevent it from reoccurring,
- subjects analyzed are not as massive compared to, for example, human health data in medicine, as risk analysis focuses on dedicated areas,
- security descriptions are usually explicitly prepared by subject matter experts (SMEs) during the dedicated analysis tasks performed regularly but rarely, specifically for selected critical infrastructure elements, i.e., power plants. Therefore, their coverage is usually limited.

On the other hand, risk prevention is a more challenging task as, for example, given the current knowledge, we would like to predict possible future risk events to prevent them *before* they happen. Hence, taking the medical domain example, given current diseases and how we treat them, we would like to predict future diseases and design their treatment *before* before anybody gets sick.

1.2.4 Validation and Explainablility

Explainable artificial intelligence (XAI) is a set of processes and methods that allows human users to comprehend and trust the results and output created by machine learning algorithms. From a risk analysis perspective, explainability is essential as a network representation of risk, which in many cases contains aggregations, must be constructed from descriptions in a verifiable way. For example, sentences: *"Astra Zeneca was first to develop a Covid19 vaccine. COVID-19 was a serious threat to global health in 2020 and* 2021" may generate a relationship between *Astra Zeneca* and *threat to global health* through *Covid19*. Therefore, a situation must be reconstructed from the original text to validate the relationship [17]. Explainability is even more critical in the case of distributed evidence, as collecting and linking them from various sources is even more challenging. From the specific NLP case perspective, for example, NLI, explainability means that not only decision: *entailment, contradiction, or neutral* is provided, but also sentences or text spans that justify it [18].

Validation is a method for evaluating the model's performance. In a basic approach, validation is performed through a dedicated data set to validate the classifier. Validation is meant to increase the model's trustworthiness and, in the risk scenario, the trustworthiness of risk-related detection. Validation is usually performed based on evaluating the classifier's performance metrics, such as Accuracy, Precision, Recall, and Sensitivity. Unfortunately, many areas, including risk analysis, lack extensive dedicated training or validation sets, and relying on 'some validation' sets designed for different cases may falsely increase the trustworthiness. The reason for this is that the meanings of words across domains differ. An example of such a case is the usage of the word "boot". It means footwear and a process of restarting electronic equipment. Validating a system on a dataset with the first "boot" meaning can harm understanding the system operating in the electronic equipment domain.

1.2.5 Goal of Reasearch

Natural language processing is a rapidly changing domain nowadays. Although the progress is exceptional, we are still far from the position to have a portfolio of components to *"plug them in"* to build a dedicated risk analysis system. Therefore, several elements are driving the research described in the presented thesis:

- First, as it is financially prohibitive, the research shall answer if it is possible to construct a risk detection system without creating a dedicated LLM for Risk Analysis.
- Second, it shall evaluate available trained and language-specialized classifiers and construct the pipeline to identify risk relations without a dedicated training set.
- Third, a validation method shall be provided.
- Finally, the pipeline shall produce the A-V-H graph, a chosen network model of risk interaction.

Overall, the pipeline shall rely on the Large Language Model's capabilities as much as possible, relaxing the problem of missing training examples and providing means to verify

the facts supporting the risk relations identified.

Chapter 2

Related Work

Drawing lessons from historical accidents and proactively identifying and even predicting the potential risks of various hazards improves the safety assurance level of any system. The thesis focuses on two aspects of this endeavor: first, consuming the risk information from natural language resources, and second, representing the risk interactions as a labeled property graph.

The system's complexity is its natural vulnerability. Normal Accident Theory (NAT) explains that some accidents are inevitable because of their complexity. The NAT approach, proposed after the nuclear accident of Three Mile Island (1979), segments systems into two broad categories: linear and non-linear [19]. In linear systems, it is possible to specifically isolate the effect and understand the impact of hazardous events. In contrast, non-linear systems are characterized by their *tight coupling* and *interactive complexity*. In the non-linear system, we encounter *normal accidents* in which multiple failures form unforeseen interactions that make accidents very difficult or impossible (with our current understanding of the system) to diagnose. Although the NAT approach has been criticized especially for having imprecise definitions and lacking criteria for quantifying system complexity [20], it marked a shift in accident analysis to "focus on the systems' properties and structure, rather than on the errors that owners, designers, and operators make in running them" [19].

In line with the NAT approach, although some of the methods were introduced before NAT, current risk analysis methods are based primarily on the decomposition of the structure of system elements. The network representation of the system models the impact of hazards and the type of applied safeguards. Depending on the degree of mathematical formalism and available information, the model may represent a general qualitative risk assessment based on, for example, the identification of causal links between failure and its effect.

The first method developed on this basis was the Failure Mode and Effect Analysis (FMEA) method) [21], which systematically inspects all possible single failure modes of individual elements of the system. Unfortunately, the quality of the FMEA model is limited

by the experience of experts and represents a de facto subjective assessment of the security system. The FMEA method is also very labor-intensive. It requires the identification of all potential events and does not allow for a comprehensive analysis of possible combinations of these hazards [5]. Therefore, the comprehensiveness and objectivity of FMEA are inevitably limited.

A systematic approach and, consequently, one that reduces subjectivity in modeling is the Hazard and Operability (HAZOP) analysis proposed in the 1960s for the chemical industry [22]. The formalism of this analysis revolves around guidewords, e.g., LESS / MORE PRESSURE / TEMPERATURE, based on which a team of experts analyses the consequences of deviation from the structurally established process flow. The identified drawbacks of this method are labor-intensiveness, a descriptive presentation of relationships between the object and hazards, and – despite the proposed formalism of the analysis – excessive reliance on the knowledge of experts. Similarly, as in the case of FMEA, labor intensiveness and the related costs harm the comprehensiveness of the model [5], [4].

In general, qualitative methods are the beginning of risk analysis, and the examples provided only illustrate the basic problems associated with them. Similar problems also exist in quantitative risk assessment models, which aim to estimate and concentrate attention on the most relevant and impactful hazards from the perspective of safety [5]. In quantitative models, the problem of identifying the system's structure, hazards, and interactions is compounded by the problem of estimating the probability of an event and its consequences. For example, the FTA approach defines a *top event* that is a critical, hazardous event (e.g., fire). Using deduction, that is, backward reasoning through a causal sequence of events over the structure of the system, FTA aims to identify basic events that lead to the top one. Apart from direct causality relation, formalism includes logical and/or gates to combine various causes over the system structure, allowing identifying the *cut set*. The cut set is a set of simultaneous *basic* events that ensure the *top* event. We can derive the *minimal cut set* that is the smallest possible (irreducible) cut set. Additionally, given the minimal cut set, the flow of the impact over the structure, and the probability distribution of basic events in the cut set, we can estimate the probability of the top event together with the *importance measures* of each basic event influencing the top event (Birnbaum and Fussell-Vesely measures) [5]. The measure will prioritize events for further analysis or safeguarding.

Although FTA represents the flow of events and provides a better understanding of potential sources of failure, it is limited. It is constructed manually, requires significant effort and knowledge, and cannot represent the interaction of several top events [5]. Ex-

cessive reliance on expert knowledge is prevalent in other quantitative methods as well. In Bayesian networks, which are more flexible than FTA as they can represent causality relations for several outcomes in the system, the flow of causality itself and conditional event probabilities or rather *believes* have to be assumed by experts [5].

Given the examples provided, it is clear that the network representation of risk interaction is not a novel approach but a foundation of more informative risk-related analysis. In parallel, as most risk-related resources are written, there is an ongoing effort to consume textual resources to comprehensively analyze unstructured or semi-structured descriptions of various risk-related events. Expert systems are already targeting this task. Although they represent risk interaction in a network manner to integrate various sources for combined analysis, they differ in framing the problem.

Formal solutions rely on representing the system, its structure, and the connectivity of its components through an ontology. This approach augments existing hazard identification processes, i.e., Failure Mode and Effect Analysis (FMEA), by defining the problem as a reasoning task. The reasoning combines facts collected through the system analysis with concepts and axioms implemented in the ontology.

Ontology is a powerful approach to detecting the impact of events on the system. This approach allows one to evaluate which components are affected by the occurrence of a specific hazard. However, it does not provide a quantitative indication of which components are the most vulnerable, meaning the largest number of hazards impacts them. The Labelled Property Graph (LPG) allows such analysis through network representation of risk interactions. The LPG allows modeling the flow of hazard in the model of the system, allowing graph algorithms, i.e., centrality measures, to indicate which nodes form, for example, hubs. Such hubs are system components that are the most vulnerable as they aggregate the impact of many hazards.

This chapter provides examples of ontology and LPG approaches to risk-interaction detection and discusses how narratives are transformed to allow the analysis.

2.1 Ontology in Risk Analysis

The definition of ontology has evolved over time. It started with one proposed by Gruber: "explicit definition of a conceptualization" [23], which emphasizes the notion of conceptualization - a structure of concepts and relations between them that is abstracted away from the real-world objects. In this aspect, conceptualization shall represent the

simplified, abstract model of the area of interest. In 1997, Borst defined an ontology as a "formal specification of a shared conceptualization" in which a "shared" feature underlines that the ontology shall be interoperable, forming a backbone of a common interpretation of the area of interest. The formalism of ontology facilitates the machine interpretability of the ontology itself, allowing automated reasoning. In 1998, Studer provided a combined definition of ontology as a "formal, explicit specification of a shared conceptualization" used today [24]. Ontology in the risk domain serves two purposes:

- to normalize (through shared conceptualization) and integrate risk information represented in various narratives and,
- to automate (as conceptualization is formal) and therefore support the impact analysis itself, as manual validation of the impact is inefficient

2.1.1 Direct Application of Reasoning on Ontologies in Risk Analysis

Application of the ontology in a formal query-based accident analysis

In the chemical industry, processes and procedures can be very complex, and descriptions of events and their contexts are difficult to interpret. However, it is important to identify cause-effect relationships and recognize lessons learned. The interpretation depends on the experience of the human experts involved in safety assessments. The ontology can help as it can be constructed explicitly around the critical analysis goals called *comptence questions* to put a narrative into a semantic framework around the goals, for example, [25]:

- "What are the hazardous events that involve a specific substance and equipment?"
- "What are potential causes of a specific hazardous event, based on the involved substance and equipment (location)?"
- "What are the potential consequences of a specific hazardous event, based on the involved substance and equipment (location)?"

Automatic ontology construction directly from the descriptions is difficult [26]. Therefore, the initial ontology can be created manually and contain the main structure of concepts and relations (a terminology box - TBOX): (Fig. 2.1):

- HazardousEvent: potentially harmful event,
- Location: involved equipment, unit or plant component,
- Substance: any involved chemical substance,
- Cause: potential causes that led to the hazardous event,







Figure 2.2: Expanded Concept Structure [25]

• Consequence: events resulting from hazardous events.

Core concepts are then used to define derived concepts, such as *Accident* (Fig. 2.2). The narrative is processed in a semi-automatic manner (as human validation and correction are needed) to link the extracted terms with specific concepts (Fig. 2.3). The goal is to put all extracted terms in relation to each other to identify cause-effect relationships: cause \rightarrow hazardous event \rightarrow consequence.

The cornerstone of such an approach is correctly identifying individuals in the text for core and expanded concepts (the assertion of terms- ABOX). In the first step, a custom pattern-based tag recognition is performed to simplify their expression for seeding concept terms. It is assumed that concepts can be described with one, two, or three words that occur in a certain sequence. For example, for the hazardous event *FIRE*, various expressions of



Figure 2.3: Creating the ontology from narratives by expanding core concepts and relations [25]

fire such as *jet fire, pool fire,* or *flash fire* will contain *FIRE* tag. For the cooling equipment, various expressions for cooling like *cooling system* or *cooling jacket* will contain *COOL* tag. The semantic relationship between terms, for example, that specific *hazardous event* took place at the *location* with specific *consequence* is assumed if they co-occur in the pre-defined context. The relationship is then manually validated.

In this way, a set of ontology individuals is created. The risk-related reasoning is performed directly on the constructed ontology with a formal query designed to answer the specific competence question and executed through the HermiT reasoner [25].

Application of Ontology to Integrate Data from Various Documents

One of the challenges in performing the risk assessment is to perform it holistically when the overall task is split into subtasks distributed across teams of experts. Failure Mode and Effect Analysis (FMEA) [27] is a good example of such a situation. From the data acquisition perspective, an established methodology constrains data to be provided to the specific FMEA format (Fig. 2.4). The format specifies that each component has its specific function within the system, its failure mode that defines how it breaks, and an associated failure effect. Unfortunately, failure analysis is performed in separate documents, and human reasoning across separate spreadsheets is laborious. Therefore, there is a need for an ontology to support it.

A hierarchical decomposition of the system into its functional sub-systems precedes the application of FMEA. Therefore, the relationships between components, which are both physical and functional, can be explicitly coded as axioms in an ontology and hence become interpretable by logical reasoner software such as OWL, Pellet, or HermiT. As each component has multiple failure modes that are related to other components, the effect

Ref.	Component	Function	Failure mode	Failure effect
20.1.1	Heaters	To heat up unit	(a) overcurrent	Loss of all heating.
			(b) short circuit	Loss of all heating.
			(c) earth fault	Loss of all heating
20.1.2	Terminal box	Connect supply to heaters	(a) overcurrent	Loss or reduction of heating
			(b) short circuit	Loss of all heating
			(c) cable failure	Loss or reduction of heating

Figure 2.4: FMEA Spreadsheet format [27]



Figure 2.5: FMEA Ontology Hierarchy [27]

of a failure at a low level can become a failure cause of a higher level. These effects are propagated through the system hierarchy until the final failure effect is identified.

In the proposed holistic analysis, the constructed ontology is the extension of ISO 15926-14 ontology [28]. It contains a target functional system decomposition (FSO Ontology), specific FMEA ontology, which represents failure effects, failure mode observations ontology (FMO Ontology), and components of a particular asset (ASO Ontology) (Fig. 2.5). The FMEA ontology defines concepts denoting the system's inferred state, for example: *ObjectInFaultState*. The reasoning is performed for a specific individual in a specific state. For example, for the object *"heater"* and the state *"heaters malfunction in the overcurrent"*. The inferred state of the *"heater"* is *ObjectInFaultState* (Fig. 2.6). Assuming the hierarchy of components within the system, the failure propagates across the hierarchy where *ObjectInFaultState* is deducted for *"heater system"* and *"heating, ventilation and cooling system"* (Fig. 2.7).

The FMEA approach is formal. It relies on formally specified ontology and the wellstructured input format. Still, the linguistic aspect of the content of the FMEA files itself blocks the general applicability [27]:

- it isn't easy to ensure that terms and relationships are used consistently, particularly when the tables are large,
- the language used is often specific to those involved in a particular FMEA activity,
- a spreadsheet typically has no explicit semantics, making it difficult to find, share or reuse the knowledge acquired during the analysis.



Figure 2.6: Sample Reasoning in the target ontology [27]



Figure 2.7: FMEA Failure propagation across the hierarchy [27]

2.1.2 Application of Ontology in the Knowledge Graph Construction

The ontology provides a shared understanding of a domain that can be used to structure information and enable interoperability between different systems and applications. The conceptualization is defined explicitly and normalizes risk interaction, which is then analyzed quantitatively through KG representation of risk instead of direct reasoning on the ontology. The ontology is then a model of the data that allows it to be stored coherently in the Knowledge Graph implemented as Labeled Property Graph (LPG).

Ontology for Railway Accident Analysis

The analysis of railway accident reports in Switzerland faces the problem of multilingual expressions of the same events or concepts (Fig. 2.9). In this scenario, the ontology normalizes the information across German, French, and Italian. Then, at the conceptual level, the risk interactions are represented in the form of a knowledge graph (Fig. 2.9) [29].

The ontology is constructed to normalize risk representation for the following *competence questions* [25] to identify incident reports in any of the source languages that relate to an injury occurring as a result of passengers:

- alighting vehicles,
- falling down stairs,
- boarding vehicles,
- being trapped by closing doors,
- being struck by falling bags.

It models the explicit relations between the document (Fig. 2.8), the *Record* node, the sentences *Sentence* node, and the words *Word* node and concepts. The *ontology* nodes represent the concept hierarchy connected by a default relation *a_type_of*, for example, *"train" is a_type_of "vehicle"* (Fig. 2.8). The relations between *terms* and *words* are confirmed manually. The analyst connects variant grammatical forms of *words* representing the same *term* through relation *forms*. The relationship *next* indicates the next word in the sentence used to identify bigram concepts. Unigram and bigram terms are identified through a measure of importance (Tf_Idf) of terms in the corpus (Eq. 2.1). Additionally, the analyst defines the hierarchy concept-term (Tab. 2.1).

The proposed approach allowed the analysis of incident reports to collectively identify the safety-related statistics around *competence questions* (Fig. 2.10). Although the proposed text processing is a rudiment dictionary of railway safety-specific terms in the three languages used in reporting, the benefits of the proposed approach are:



Figure 2.8: Multilingual Ontology with document structure[29]

- · significantly simplified query design executed on the KB instead of ontology,
- quantitative analysis of cases of injuries.

$$TF_IDF = \frac{f_T}{N} ln(\frac{R}{R_T})$$
(2.1)

where:

- f_T is the frequency of occurrence of term T within a record,
- N is the total number of terms in the corpus,
- R is the number of records in the corpus,
- R_T is the number of records that contain the term T.

Ontology and Knowledge Graph Construction for Near-Misses

Descriptions of "near-misses" provide detailed information on complex hazardous situations, including their initialization, evolution in the system, and barrier effectiveness. This data is critical from a safety analysis perspective, as barriers may eventually fail in similar situations.

To properly extract risk interactions from narratives, the use case for ontology is to explicitly model the event's evolution in the system and the effectiveness of the barriers [7]. The goal of the KG representation is to answer the following competence questions:



Figure 2.9: Multilinqual Concept Expression [29]

Concept	Term
vehicle	carriage, vehicle, ambulance, tram, train, bus
person	doctor, self, customer, person, driver, passenger, months old, years old, baby, young, old, female, male
object	Bag, alcohol, drugs, stairs, footboard, customer information system, ticket, door

Table 2.1: Concept - Term mapping [29]

Query	Language			
	German	French	Italian	
1. Alighting	693	296	34	
2. Falling down stairs	73	9	1	
3. Boarding	1100	349	16	
4. Closing doors	1220	464	409	
5. Falling bags	3	1	0	

Figure 2.10: Multi-linqual railway incident incident results [29]


Figure 2.11: EsOpAI Ontology Structure [30]

- "Which are the more vulnerable apparatus in the system (e.g., refinery)?
- "How do they usually fail?"
- "Which are the consequences of such failures?".
- "Which is the most critical barrier preventing the catastrophe?".

The ontology (Fig. 2.11) is constructed using concepts defined in EsOpAI (Operational Experiences via Artificial Intelligence) [30] to represent a near-miss. It contains the following entities: EVENT, APPARATUS, SUBSTANCE, ACTIVITY, BARRIER, and PEOPLE. The concepts are further specified to represent more granular concepts. For example, EVENT is split into:

- LOSS related to losses of containment (e.g., leakage, overfilling), FAILURE, which contains both failure and damages,
- DETERIORATION that includes all the mechanisms of deterioration (e.g., corrosion, pitting, creep) that caused integrity problems,
- MAJOR depicting all those events which have the potential to generate other incidents (e.g., fire)
- SUCCESS to indicate the positive barrier action that contributes to interrupting the incident escalation

The APPARATUS entity is divided into EQUIPMENT, which indicates a whole, and COM-PONENT, which indicates a part of the EQUIPMENT connected with a PART_OF relation. EsOpAI has been designed with four relationships:

• RELATED_TO, which is a generic relation between two entities,

Relation	Head Concept	Tail Concept
causes	EVENT, ACTIVITY, PEOPLE	EVENT
involves	EVENT, ACTIVITY, APPARATUS, BARRIER	SUBSTANCE
part_of	APPARATUS, BARRIER	APPARATUS, BARRIER
related_to	EVENT, BARRIER, ACTIVITY	ACTIVITY, PEOPLE, APPARATUS, BARRIER

Table 2.2: EsOpAI Relations [30]



Figure 2.12: EsOpAI Annotation Example [30]

- PART_OF describing a physical connection between two entities,
- INVOLVES, that relates an entity to a substance,
- CAUSES, which states a causal connection between two entities

The relations are predefined and established between the entities (Tab 2.2). Supervised learning is used to detect concepts and relations in the narrative. The annotation strategy strictly follows the structure of the ontology, and the entities and relations are annotated exactly to the intended concepts and relations (Fig. 2.12). Finally, detected triples i.e.: *head_concept - relation - tail_concept* are extracted directly into the KG. The structure of the KG is aligned with the ontology as the triplets loaded are detected according to the ontology definition of entities and relations between them. The KG allows for "ontological explorative analysis" [7], a combined network analysis of numerous near-misses across installations and plants. It allows quantitative answers to the competence questions, for example, "most frequent causes of failures upon the most frequent apparatus that fail in refineries" (Fig. 2.13).

AE	f_{AE}	$oldsymbol{f}^{\%}{}_{AE}$
FLANGE LOSS	7	3.30%
PIPELINE LOSS	7	3.30%
GASKET BREAKAGE	6	2.90%
VALVE LOSS	4	1.90%
PIPE DAMAGE	3	1.40%

Figure 2.13: Top five apparatus failures (AE) causes in refineries in Italy [7]

2.2 Knowledge Graphs and Network Representation of Risk

The network representation of risk enables the topological analysis of its interactions: degrees, shortest paths, and centralities [1]. For example, the betweenness centrality measure will indicate the degree to which the node transmits the impact of the hazard to other nodes. Identifying such nodes in the model would help focus the protection measures around that specific system element. In general, network representation of risks will allow [1]:

- to characterize the risk-related impact of each system component, represented as a node, in the system's overall structure,
- to simulate the impact of the system's structural changes from the risk control perspective. For example, the analysis of the impact of change in strictly law-regulated IT infrastructure due to the statutory or standards change of the selected component.
- to build a simulator of real situations that are not repetitive (crisis situations). The simulation can occur under varying conditions of risk.

There are several examples of applications of such risk representation strategy. They differ in the method for normalization of risk relations between the intended concepts. These methods do not rely on prior ontology definitions. The assumed KG structure defines the Knowledge Acquisition Pipeline, which consumes textual information directly and transforms it into a required representation. Entities and relations between them are key elements to detect.

2.2.1 Representing HAZOP Safety Reports as Knowledge Graphs

Significant safety knowledge is encoded in textual sources, such as formalized safety reports, such as the HAZOP, near-misses or FMEA. However, this knowledge remains undiscovered because it is not transformed into a representation suitable for comprehensive

analysis. The Knowledge Graph (KG) representation allows such analysis. The structure of the KG depends on the modeling goals, namely, which aspect of the narrative is useful and shall be represented. For example, the risk interaction network (Fig. 2.14) consists of four node types: suggestion, result, cause, and equipment/assets. The goal of the representation is *evaluate and combine* hazard causes, effects, and prevention measures and, therefore, validate the gaps in the safety design. Apart from that, in this specific case, the graph allows employees to picture and understand the safety and the operational requirements for the process [31].

The acquisition pipeline assumes a specific input data format i.e. HAZOP. HAZOP, which is a semi-structured analysis methodology, uses a defined set of guidewords to help evaluate the consequences of deviation from the regular process flow [5] simplifying entity detection [31]. In this case, processing is performed as follows. There are three application layers. The conceptual layer encodes knowledge of the processes used in the extraction performed in the extraction layer. It aims to normalize HAZOP reports to combine information from different processes and represent industrial safety knowledge. The goal of the layer is to define the node types (concepts) used in the storage layer, like cause, result, equipment, or suggestion. For example, a cause entity could be: "coil blockage", result: "crude oil in the vaporizer will boil" and the recommended entity: "clean the coils". The extraction layer extracts all entities and relationships among them. In short, it extracts the knowledge triple: subject, relationship, and object and saves them into a storage layer for reporting. The storage layer is the LPG containing the network representation of risk (Fig. 2.14).

2.2.2 Representing Free Text Risk-Related Narratives as Knowledge Graphs

Reusing and releasing the value of industrial safety knowledge is a step towards a comprehensive risk repository. Modeling risk interaction as a heterogeneous graph, where the node type represents a concept of risk, cause, or asset, increases the clarity of representation. However, the approach proposed for the HAZOP analysis is limited because it assumes a specific input format; hence, the processing has to be aligned.

A solution detecting risk relations for railway safety parses natural text description of British rail accidents [6]. The processing is step-wise (called the knowledge extraction steps), where each step is responsible for detecting, linking, normalizing representation across documents and mentions (knowledge fusion step), and eventually representing risk interactions in the form of the knowledge graph, namely RKGRS (Fig. 2.15).

The resulting structure represents risk interactions as a heterogeneous graph with



Figure 2.14: Exemplary Storage Layer Content for Hazop Analysis [31]

defined node types representing different concepts. For example, C nodes indicate Cause, D - danger, and K consequences. For example, C01 (imperfect management of maintenance practices) will result in danger nodes: D01: struck by object or collision, D07: derailment. Consequences are K03: damage to structure, component, or device, K01: injuries (Fig. 2.16).

The ability to parse free text reports is a step in the right direction, as the input format shall not limit the comprehensiveness of the risk repository. However, relaxing the format results in a complicated knowledge acquisition pipeline. In this case, the ensemble of classifiers is trained in a supervised manner [6], and the training set was designed specifically for the railway safety scenario. The training uses **seventeen** text augmentation algorithms to enrich the representation (through embeddings) or disturb the text insignificantly by adding or replacing characters or synonyms to achieve acceptable training results. To apply the proposed approach to other domains, the training set and text augmentations shall be adjusted [6].

2.3 Summary

The methodology for integrating and representing risk interactions from description is not established yet. Solutions are specific and target specific areas, such as railway accidents



Figure 2.15: Railway Accident Reports Processing Pipeline [6]



Figure 2.16: Railway Safety Knowledge Graph [6]

or specific competence questions for the process industry. The limited applicability of current solutions does not come from the fact that creating a "comprehensive risk ontology" that would normalize risk representation across various descriptions is potentially impossible. The ontology is a *shared* conceptualization. Therefore, with a coordinated effort, we can potentially cover the most critical, targeted areas such as chemical processes, railway, or electricity generation. Within the "comprehensive ontology," it would be possible to run simulations to identify areas that require specific protection.

The flexibility of human language in risk description is a limiting factor. This flexibility translates into complicated text processing pipelines, and it lacks a general, comprehensive approach to detecting risk-related relations in the narrative. In many cases, domain-specific solutions are required to prepare their own detection mechanism, which relies on the specific training and validation sets. Even solutions targeting well-defined areas need to address the training set issue by providing additional constraints on data, i.e., hierarchy [32]. This approach uses the taxonomy of events to augment the training and improve the detection quality. Examples of an element lower in the hierarchy will be used to train higher-level concepts.

Focusing on the linguistic aspect of risk modeling is necessary and may become a foundation for an application that comprehensively represents risk interaction.

Chapter 3

Relevant Natural Language Processing Techniques

The examples described in the previous chapter relied on two key natural language processing elements: Entity Recognition (ER), which is responsible for detecting a concept such as "event," "apparatus," "cause," or "suggestion" in the narrative, and Relationship Extraction (RE), which confirms risk-related semantic relations between detected entities. This section describes NLP techniques essential in the ER and RE areas. It discusses their limitations and applicability in detecting risk-related entities and relations.

3.1 Entity Recognition

Entity Recognition, or Named Entity Recognition (NER), is an NLP task that aims to assign a text span to one of the pre-defined semantic categories like 'PERSON', 'COUNTRY', 'ORGANIZATION' [33]. NER is useful in applications requiring such classification. In KG construction, NER assigns the concept and classifies nodes. For example, the knowledge acquisition pipeline identifies "Danger", "Consequence", and "Cause" categories to represent concepts in railway risk interaction. However, detecting risk-related categories is difficult to obtain using current NER approaches.

Formally, the NER classification is the sequential prediction problem to estimate the probabilities of predicting the ith BIO tag given the *context* being the history of *k*, the future

of *l* words, and the history of *m* past BIO tags:

$$P(y_i|x_{(i-k)},...,x_{(i+l)},y_{(i-m)},...,y_{(i-1)})$$

where k, l, and m are small numbers. Algorithms used to estimate the probabilities of the tags are:

- Conditional Random Fields (CRFs): CRFs model the conditional probability of label sequences given input features, capturing dependencies between neighboring labels,
- Bidirectional LSTM (BiLSTM): BiLSTM networks are recurrent neural networks that can capture sequential dependencies in the input data,
- Transformer-Based Models: Transformer architectures such as BERT [35], GPT [13], and their variants which have shown state-of-the-art performance in NER by leveraging self-attention mechanisms to capture contextual information

However, to design an efficient NER detection system, the following design questions need to be answered [36]:

- · How to model non-local dependencies,
- How to use external knowledge resources or construct a training set.

From the RA perspective, both pose significant difficulties.

The non-local dependencies mean that the classification of the text span depends on the selection of the context and changes with the context size. The local dependency means that syntactic and semantic features of the entity must be present in the current context of k past l future words and m past tags. Such a context definition is fixed across the narrative to perform the classification. For example, in the sentence "Joseph Robinette Biden is the 48th President of the United States of America." term "America" is a part of the *LOCATION*, and all features required to classify "America" this way are present in the context. Likewise, in the sentence "Brian Moynihan is the President of the Bank of America", the term "America" is a part *CORPORATION* as local features of the same context definition allow classifying the term correctly though the different class is assigned.

The fixed, local context assumption harms the classification in the risk-analysis domain. In the railway case [6], to classify a text span to a "Danger" category, the context shall contain features expressing some harmful effect of it, as the token "Bank" differentiates the classification of *America* in previous examples. In general, such features may not follow the locality assumption as effect, e.g., injuries and damage to the infrastructure, can follow the "derailment" danger in an arbitrary long context.

Another difficulty in a direct application of the Entity Recognition approach is the *contextuality* of risk-related concepts. The *contextuality* means that the classification of the

text span depends on the classification of other text spans in the context. Formally, current tag classification depends on the fixed history of *m* tags. Depending on the intended focus, the same text span can have multiple classifications in the same context. For example, in the sentence "A spark ignited a container supplying fuel to the engine," depending on the focus, the "container" is a "Danger" as it can be ignited by the spark, potentially destroying the engine. Simoultenousely, it is a "Component" as, without it, proper engine functioning is impossible as it is responsible for fuel supply.

Last but not least, supervised learning is the dominant approach to solving the classification task. In the risk analysis domain, standardized training sets are not available. Preparing a training set means manually solving the entity recognition task, especially in rare, extraordinary malfunction cases.

3.2 Relationship Extraction

The efficient, algorithmic Relationship Extraction started with introducing the Hearts Patterns. Hearst proposed to detect hyponymy, namely "is-a", relations, e.g., "rose is a flower". The approach used a set of regular expressions on the syntactical decomposition of the sentence. It was motivated by two goals:

- · to avoid the need to pre-encode the extensive knowledge and
- to apply the same approach across a wide range of text.

[37]. The key assumption was that the syntactical patterns are enough to define relations within a sentence effectively. It quickly turned out that the flexibility of human language was underestimated, and the pattern approach proved ineffective. The Hearst Patterns, however, initiated intensive research in the RE area. This section will evaluate the applicability of existing RE methods in the context of scarce resources in risk analysis.

3.2.1 Intra-Sentence Relationship Extraction

Sentences are composed of smaller linguistic units, such as words and phrases, and the meaning of a sentence is determined by the meanings of its constituent parts and the way they are combined. Sentences convey semantic relationships between their parts, words, and phrases that can be detected and identified. Their syntactical decomposition provides additional features: parts of speech, grammatical relations, or named entities within the sentence (Fig. 3.2) [38]. These features were used extensively in the initial approaches to the RE.



Figure 3.1: Exemplary Sentence



Figure 3.2: A sample sentence decomposition to a feature-rich dependency tree

Rapid progress in the RE started after launching a dedicated relationship extraction task between pairs of nominal within a single sentence [39]. Initial solutions relied on heavily syntactic features and formulated the problem as kernel-based classification of the shortest path in the sentence dependency tree [40], maximum entropy models over syntactic features in the sentence [41] or graphical models [42]. Currently, two approaches, a transformer-based [43], and dependency decomposition, graph-based [44] trained in a supervised manner, are currently the best models.

The main drawback of the state-of-the-art approaches that limit their direct applicability in the Risk Analysis domain comes directly from the training strategy. A supervised approach, although very effective [43], [44], does not guarantee the same performance on the vocabulary that is out of domain [18]. It is hard to imagine a scenario in which each new failure event is provided with a dedicated training set that would contain annotated elements. This would mean we solve RE classification tasks manually each time we analyze risk interactions for a specific case.

Relaxing the supervised training strategy and relying on an unsupervised approach is tempting. However, the unsupervised approach falls short of applicability to risk analysis due to large corpus requirements or reliance on auxiliary classifications, such as entity



Figure 3.3: Verbalizer Approach to RE [48]

recognition, to improve clustering results [45].

An approach that relaxes the training set and the corpus requirements uses the direct application of textual entailment [46], [47] in relation detection. In this approach, a relation classification task is cast as an inference of the hypothesis that the sentence, being a premise, entails the relation pattern of interest [48]. The approach assumes that the subject and the object of the relation are in the sentence, and relationship templates are augmented to validate them against the premise. A template with the highest entailment score is selected (Fig. 3.3) [48].

Although relaxing the most critical constraints on risk analysis, namely training set and corpus, the verbalizer is limited in the assumption that the sentence defines a relationship and the patterns are evaluated for sentences only. In a general corpus, almost 40% of relations are defined across sentences [49]. Unfortunately, the distributed nature of the relations in the risk domain is prevalent as information on hazards is stored across documents [31], [6].

3.2.2 Inter-Sentence Relationship Extraction

Inter-sentence relationship extraction is an approach that focuses on identifying relations across the whole document. Such detection requires complex cross-sentence inferences to synthesize information from the whole document to detect relationships correctly. Such reasoning is relevant in the domain of risk analysis as well and includes [49]:

- logical reasoning, that requires identification of proper *"bridge entities"* that connect the head and tail of the relation. Such specific entities are essential to model the propagation of hazards in the description of the system.
- coreference reasoning, that requires correct identification of head and tail entities in



Figure 3.4: Inter-sentence Relationship Example

the whole narrative,

• common-sense reasoning that requires some additional knowledge to identify a relation. For example, "William and Kate had 4 children". Therefore, Kate could be *"a spouse of"* William. In the risk domain, such scenarios are not uncommon. For example, *"gasoline can ignite"*. Therefore, assuming that gasoline *"can cause"* fire is natural.

Therefore, identifying relations in an inter-sentence scenario requires modeling semantic interactions between mentions of entities between sentences. The direction of the relation may not follow a "reading direction." For example, it is possible, given the exemplary text (Fig. 3.4.), to verify that "droplet" is in a relationship with the "engine". In this case, the analysis may follow the deduction path: "droplet" \rightarrow "supercooled water droplets collide with a surface" \rightarrow "if supercooled water droplets collide with a surface they may result in blocked fuel inlet pipes" \rightarrow "blocked fuel inlet pipes" \rightarrow "aviation fuel designed for use in aircraft powered by gas turbine engines" \rightarrow "engine" Document-level relationship extraction will be reviewed for two types of models: sequence-based and graph-based.

Sequence-based models

Sequence-based approaches avoid explicit modeling of the document. Instead, a deep neural structure is created, where each component specializes in a dedicated inference separating the relationship within sentence (local) and document (global) (Fig 3.5) relationship detection. The specific architectures use different neural networks, i.e., Bi-LSTM, to model local and global dependencies between entities [50]. The document-level Bi-LSTM layer is responsible for modeling the multi-hop, cross-sentence reasoning.

The separation of contexts models interactions document-wise and supports crosssentence reasoning. The same philosophy, but using another type of neural networks approach, has been proposed in the Mentioned-based Reasoning Network (MRN) [51]. In this approach, local contexts capturing close subject-object relations are modeled through the convolution of entity mentions. Stacking the convolution models multi-hop dependencies between them in the narrative. A document-wise representation of a relation, a global



Figure 3.5: Architecture of the Hierarchical Inference Network. Sentence-level, entity-level context is separated from document-level context [50]

context, is achieved through a co-attention of local context convolutions. Co-attention performs a weighted combination of several local contexts (local relations) to achieve its global representation.

Another approach in sequence models relies on the observation that most crosssentence relations are fully defined within the fixed context anyway. The statistics around the general distribution of head and tail relation entities is that they are mostly separated by three sentences in the narrative [52]. Therefore, the RE detection is cast as selecting a combination of at most three sentences from the document. These sentences would form the relation context. The context must include references to both the head and tail entities, and it is used to classify the relationship. The strategies for selecting the sentences are as follows. Sentences may form a "Consecutive Path" (Fig. (3.6) following each other in the "reading order." Another type of context supports multi-hop reasoning, meaning each sentence is connected by a common entity. In this context, sentences may not follow the "reading order"; however, their number is still limited to three (Fig. (3.6). The third strategy defines the context as pairs of sentences containing the first, the head entity, and the second, the tail. The set of contexts is the cartesian product of sentences containing the entities (Fig. (3.6). The contexts constructed are used to train a discriminative classifier based on the prepared training set [52].

The sequential models are discriminative classifiers trained on a dedicated documentlevel training set.



Figure 3.6: Types of contexts in sequence-based document level RE [52]

Graph-based models

A graph-based approach has gained significant attention lately and is considered the most effective approach to document-level RE. It is based on a network representation of the document [53]. The representation captures semantic, syntactic, and positional information on the entities, providing more features for relationship classification. The general processing pipeline involves splitting documents into sentences and **detecting the entities** that will be evaluated for relationships in the document. The document-level relationship classification tasks are usually cast as a link prediction problem between the nodes representing the entities of interest.

There is a significant number of models solving document-level RE tasks varying in implementation details. For example, the Edge-oriented Graph (EoG), which is the extension of an earlier model for solving sentence level RE [54], represents the document as a heterogenous, undirected graph containing the following node types (Fig. 3.7):

- Entity node represents concepts.
- Mention node is a span describing the entity.
- Sentence node is a sentence that holds the mentions.

Each node representation is computed as the average of the embedding of its constituent words. The embeddings of words are computed in the Sentence Encoding Layer using the BI-LSTM network (Fig. 3.8). The edges represent different roles of the node in the document:

- Entity-Mention (EM) edge indicates that the mention describes an entity,
- Mention-Sentence (SS) edge indicates that specific mention is a part of the sentence,
- Entity-Sentence (ES) edge indicates that the entity is contained in the sentence
- Sentence-Sentence (SS) edge connects the sentence to model non-local dependencies. The edge specifies how many other sentences (the distance) separate the sen-







Figure 3.8: EoG Architecture [54]

tences in the document.

• Mention-Mention (MM) edge connects mentions that are a part of the same sentence. . The edge representation is calculated as the concatenation of the representation of connected nodes. The link prediction task is cast as the classification of the path that links entity nodes in the network. The path representation is calculated as the non-linear transformation of embeddings of edges in the path [54].

In the Global-To-Local Neural Network for Relationship Extraction (GLRE) [55], the document structure is also explicitly represented as a network (Fig. 3.9). The structure of the graph remains similar to EoG's with only slight modification in that all sentences are connected without addition *distance* information:

- Entity-Mention (EM) edge indicates that the mention describes an entity,
- Mention-Sentence (SS) edge indicates that specific mention is a part of the sentence,



Figure 3.9: GLRE Document Representation and Classification Schema



Figure 3.10: Global-to-Local Neural Network for RE Architecture [55]

- Entity-Sentence (ES) edge indicates that the entity is contained in the sentence,
- Sentence-Sentence (SS) edge' connects all sentences regardless of their distance,

• Mention-Mention (MM) edge connects mentions that are a part of the same sentence. The representation of nodes is an average of the embeddings of its constituent words. Words embeddings are calculated during text transformation, and it is the embedding calculated through the BERT transformer of the sentence or short text fragment for better contextuality [35] [55]. There is no representation associated with the edges. The relation detection problem is cast as a link prediction problem between entity nodes.

There are two representations of entities: global and local. Each entity can have only one global representation that is convoluted with the network structure of the document [56]. Local representations are weighted combinations of local and global embeddings of their mentions (Fig. 3.9). The dedicated multi-head attention modules (Multi-head Attention 0 for the head and multi-head Attention 1 for the tail) combine the mentions. The final link prediction task between head and tail entities uses concatenated embeddings of both in the logistic regression task. Contrary to the previous example, which used a representation of the path between the nodes, the GLRE approach combines document structure and local linguistic features of the mentions in the RE classification task. [55].

The efficacy of current document-level extraction is limited specifically in two ways:

- it is supervised and may not be generalized onto documents semantically far from the training sets.
- the approach relies on the pre-processing which detects *entities* and *"entities span"*. Quality of entity detection is essential for the overall performance [57].

3.2.3 Transformer-based Relationship Extraction

In LM domain, the RE classification task can be cast as another NLP task. For example, the verbalizer casts RE as an entailment task relying on the generalization capabilities of LM in entailment recognition [48]. The other one, but chronologically the first, frames RE classification as a question-answering (QA) task [58]. This is the first approach to a zero-shot classification scenario in which a classifier is used for cases not presented in the training phase. The leading assumption of the approach is that question formulation generalizes the structure of the relation. Therefore, training on question-answer pairs instead of relation examples is more efficient. In the direct RE training, the classifier is responsible for identifying linguistic elements defining the relation abstracting away from combinations of subject and object [43]. In the QA scenario, the classifier is trained on an extensive QA dataset instead of an RE dataset [58].

The RE classification problem is cast as a parameterized, relation-specific question, and the classifier identifies span as an answer. An answer is an object of the relation (Fig. 3.11). The relation-specific question's parameterization must be specified (Tab. 3.1). For example, given the sentence *Brian Moynihan is the President of the Bank of America*, the goal would be to evaluate the workplace for *Brian Moynihan*. The parametrized question would be *Where does Brian Moynihan work*. In case of a relation not supported by the text, the classifier returns an empty span.

Language Models are trained extensively over a large corpus. While learning, some relational knowledge is encoded and can be retrieved directly through a masking strategy.



Figure 3.11: Question Answering Architecture for BERT [35]

Relation	Question Template
educated_at(x,y)	Where did x graduate from? In which university did x study?
occupation_at(x,y)	What does x do? What is x's job?

Table 3.1: Question-Relation Templates [58]



Figure 3.12: Language Models as Knowledge Encoders [59]

In this scenario, the LM is considered as a linguistic memory [59]. To access this memory, a simple relation pattern is evaluated. For example, let's assume that the fact that *Brian Moynihan* works for *the Bank of America* and this information has been encoded in the LM during its training. If a templated query using a masking strategy is formulated to query the workplace: *Brian Moynihan works for [MASK]*, then the LM shall substitute *bank* as the highest probable token associated with *[MASK]* (Fig. 3.12).

Large Langauge Models, i.e., GPT3, also encode language knowledge. However, the conditional, autoregressive text generation capabilities [13] impact the approach to RE. The conditional, autoregressive generation means that response is generated given the buffer. The buffer may contain the description evaluated for relations, some instructions, and examples, which is **the prompt** [13] [60]. Conditioning on the buffer, information on the relations will be in the generated output. Specifically, the prompt opens a new avenue to provide either relationship examples or instructions on how relations shall be detected in the provided text.

A naive prompt engineering approach (Fig. 3.13) for RE initially did not yield reasonable results, mostly because the prompts were not addressing the language phenomena, which led the classifier to focus on shallow language features, which could be, for example, a simple overlapping of words. The reason for performance significantly below fine-tuned LM (BERT) were identified to be:

• low relevance regarding entity and relation in the existing sentence-level demonstration, Given the possible relations: [member of, field of work, work location, ..., father, sibling]. What are the relations between the subject entity and the object entity expressed by the sentence? **Sentence:** Savi was born in Pisa, son of Gaetano Savi, professor of Botany at the University of Pisa. **Subject:** Gaetano Savi **Object:** Botany **Relation:** field of work

Figure 3.13: Naive Prompting [62]

• the lack of explaining the exemplary mappings of demonstrations via precise instructions in natural language.

[61].

A significant improvement in LLM RE has been achieved by designing prompts combining instructions and examples for the target relation [62]. Therefore RE task has been formulated as a prompt-generation task containing:

- head and tail entities of the relation,
- examples of the relations from the repository of relation mentions,
- instruction to retrieve each relation in the provided text according to the examples or "no relation" in case nothing can be matched.

There are several drawbacks to the direct application of LLMs in RE:

- in-context approach itself. Although significant progress has been made in context size, it is still limited. Additionally, the bigger the context, the more challenging the RE task. There is no additional information on the structure of the document. Compared to the graph RE approach, the network representation of the document defines the contextuality and connectivity of entity nodes. Although the network structure may be considered an *inductive bias* of the graph approach, it provides additional information that is missing if the whole document is provided as a part of the prompt. In this case, the LLM itself must select the fragments relevant from the perspective of the specific relationship.
- LLMs will not provide a measure indicating the detection quality.

3.3 Textual Entailment

Inference or entailment is a critical ability to draw conclusions. Ido Dagan defines textual entailment as a directional relationship between pairs of text expressions, denoted by T (the entailing "Text") and H (the entailed "Hypothesis"). It is considered that T entails H if *a human reading* of T typically infers that H is most likely true [46] [63]. In the past, there used to be several domain-specific applications that were running language inference like textual entailment, and Textual Entailment (TE) is an attempt to provide a generic framework that would define the mechanisms of such semantic inference across many domains and establish a coherent evaluation of the proposed inference mechanisms [46].

In essence, textual entailment is a relaxation of the formal logical entailment and comprises several elements that come directly from natural language, or *human* perception of it, namely, "what a person would typically infer from the premise.". The textual entailment traits can be put into the following categories:

- generalization, which represents the hypothesis as a more general statement of the premise, e.g., premise: *"antibiotics inhibit the synthesis of bacterial cell walls."*, hypothesis: *"antibiotics slow down the development of bacterias"*
- inference, which involves deriving *new* facts and grounding them with the provided premise using, e.g., logical reasoning (premise: *"dealer sold 103 cars"*, hypothesis: *"dealer sold over 100 cars"*).
- paraphrasing is the situation in which premise and hypothesis are equivalent, and the textual entailment is a bi-directional relationship. For example, premise: "Frosty situations lead to conflicts", hypothesis: "Unfriendly situations lead to tensions or animosities"

The textual entailment is a three-way classification task in which the system shall detect if a given premise/hypothesis pair is either *entailed*, *contradicted*, or *neutral*. The exact formal algorithm that would assign the pair to either of the classes does not yet exist. From the algorithmic point of view, entailment identification is a complex "NLP complete" task [63] that involves multiple techniques to run the inference. For example, to identify an entailment for hypothesis 1 (Fig. 3.14), it is required to run the inference to confirm that "BMI" acquired an "American concern". It requires reasoning that "Huston" is the capital of Texas and "Texas" is a part of the "USA", a synonym of "America". However, this is not enough. The headquarters location in Houston still does not make the company American. To complete the inference, several additional connections must be established: between the LexCorp owners and the fact that the Americans live in Huston, and that concern is a coreference to LexCorp itself.

A series of workshops drove the early attempts to solve the TE problem, "The PAS-CAL Recognising Textual Entailment Challenge." Initially, solutions relied on the syntactic **Text:** The purchase of Houston-based LexCorp by BMI for \$2Bn prompted widespread sell-offs by traders as they sought to minimize exposure. LexCorp had been an employee-owned concern since 2008.

Hyp 1: BMI acquired an American company.Hyp 2: BMI bought employee-owned LexCorp for \$3.4Bn.Hyp 3: BMI is an employee-owned concern.

Figure 3.14: Entailment Example [63]



Figure 3.15: Alignment of similar terms for T (top) and H (bottom) [63]

decomposition of premise and hypothesis to match terms, called *anchors*. Generally, the process was structured as follows [63]:

- Candidate Alignment Generation. In this step, premise and hypothesis texts are chunked to select pairs of terms from both (Fig. 3.15). For each pair, the similarity score was calculated. The similarity score could be as simple as 1 for identical terms and 0 otherwise. However, for example, the candidate (1), "purchase", "acquired" are synonyms (Fig. 3.15). There is no candidate "purchase", "BMI" as the terms are of different parts of speech. Therefore, in reality, more sophisticated similarity functions were used. [63].
- Alignment. This step selects the best alignment between terms in premise and hypothesis through, for example, greedy maximization of similarity score.
- Classification. Given the aligned anchors, a feature vector is constructed. The feature vector represents "the state of similarity" between the premise and hypothesis and is used in supervised training. Given the training set, the classifier learns the decision on entailment / non-entailment given the representation of the premise and hypothesis.

Early attempts to solve the TE problem relied heavily on feature engineering on the lexical and syntactical decomposition of T and H, therefore their generalization capabilities were limited. Transformer [9] approach casts the TE task as a fine-tuning task on top of the pre-trained Language Model (BERT or RoBERT) [35], [64]. Throughout the pre-training

phase, the transformer model acquires significant linguistic abilities that relax the extensive manual feature engineering requirement. These abilities are:

- 1. effective and thorough representations for the meanings of sentences (i.e., their lexical and compositional semantics) which can be observed through contextualized embeddings of words [35]. It means that embedding of the word "bank" would be different in the context of "river" and "money".
- 2. ability to handle lexical entailment. For example, the relation between "cat" and "animal".
- 3. ability to handle quantification, as the transformer is capable of discerning the similarity of sentences based on words such as "none," "some," and "few." These differences will be visible in the sentence embedding.
- 4. and much more, for example, the ability to handle coreference (through attention mechanism), differentiate tense expressions as it can distinguish past, present, and future expressions, modality (which expresses the possibility the statement is true i.e. would, could, or possibly), and lexical and syntactic ambiguity which means that two sentences expressing same thing but differently will have embeddings close to each other.

These abilities are encoded into embedding the 'CLS' token, which is the output of the Transformer component of the classifier (Fig. 3.16).

In addition to extensive pre-training that captures linguistic phenomena, LM is finetuned to the RTE task. The fine-tuning is performed in a supervised manner (Fig. 3.16), with two large training sets: Stanford Natural Language Inference (SNLI) dataset [47], and The Multi-Genre NLI Corpus (MNLI) [65]. These datasets address the limitations of earlier RTE training sets, such as:

- they were limited in number of training examples. Early RTE corpuses had a couple of thousand hand-labeled examples. This is not enough for deep neural and transformer models.
- from the linguistic perspective, the examples were not diverse enough and, in many cases, were simply incorrect. For example, early examples referred to coreference resolution, which was problematic to resolve even for annotators [47]. An example of such a case that assumes a specific interpretation of "New York" as a city instead of the state would be the premise: "A tourist visited New York." and the hypothesis: "A tourist visited a city."



Figure 3.16: Transformer Architecture for RTE with training

• MNLI specifically provides entailment examples from ten different sources across genres and domains to model the usage of modern American English [65]. It significantly improved the generalization capabilities of the classifier as the training set accounted for more diverse examples of entailment expressions.

3.4 Semantic Frames And Semantic Role Labeling

Charles J. Fillmore originally introduced Frame Semantics with the basic idea that one cannot understand the meaning of a single word without access to all the essential knowledge related to that word, namely, its semantic frame [66]. The semantic frames are strictly associated with the cases expressed in the sentence, and at times, the innovation was that it called for the case organization of sentences, known as the *compositionality assumption*. In other words, sentences consist of cases that define specific roles of nouns, e.g., Patient, Agent, or Instrument, that support cases. Cases were considered events or scenes to study the semantics of the words involved. The approach has led to the creation of FrameNet. This repository is the evidence of the semantic and syntactic structure of the cases considered semantic building blocks for each sentence.

In the FrameNet, each Semantic Frame was defined as the specific case and the list of its arguments. The Frame consists of Frame Elements (FE) and the Lexical Units (LU). The Frame Elements are participants (case roles) defining the frame, and LUs are text fragments that evoke a given frame. For example, the frame *commercial transaction* would be evoked



Figure 3.17: General Semantic Roles

by the LU John is buying a new car from 20.000 USD, and its FEs would be:

- Buyer or Agent: The person or entity making the purchase, John in our example,
- *Seller or Patient*: The person or entity selling the product or service, in our example undefined,
- Product/Service or Instrument: What is being bought, car
- Price: The amount of money exchanged for the product or service, 20000USD.

The same frame *commercial transaction* would be evoked for synonyms of *buy*, i.e., *purchase, acquire* etc. The FEs define "Who, What, Where, When, With What, Why, How" for each frame. The precise semantic detection of frame elements (Buyer, Seller, Price etc), in general, is not possible. For example, in the sentences *I ate dinner with Anna* and *I ate dinner with sticks*, the objects defined with "with" are difficult to distinguish for their specific roles. Therefore, generalized roles, e.g., ARG0, ARG1, ARG2, and others, are currently detected [67] (Fig. 3.17); however, the frame structure is preserved, and the assumption of case compositionality of the sentences holds.

The compositionality assumption means that each sentence is the combination of its frames. Therefore, a semantic relation between the entities can be performed within the frame's arguments after frames, and their arguments are extracted from the sentence.

Summarizing, a semantic frame serves several purposes:

- it is a means to abstract cognitive schemata, and it is this schema computational counterpart [68], which means that each frame encodes the relationship and its arguments,
- The structure of the semantic frame naturally identifies its subject and objects (as they are the frame's arguments) and the relationship between them [17] [69], which simplifies the detection of semantic relations,
- The structure of semantic frames within a sentence decomposes the semantics of the sentence. Therefore, the relationship detection can be performed directly on

elements of the frames, not the cartesian product of elements of the sentence itself.

Semantic Role Labelling (SRL) [70] [71] identifies the frame's structure. A state-of-theart deep pre-trained SRL classifier [72] detects the simplified structure of a frame where instead of an agent, a patient, or an instrument, it detects generic simplified arguments of a verb: ARG0, ARG1, and others [67].

3.5 Summary

In the Natural Language Processing domain, it is a known fact that solutions explicitly designed to model specific scenarios are not guaranteed to be generalized in another case as they rely on the detection schema (classifier and the training set) targeted for that case only. Although it might sound reasonable to apply more sophisticated language representation as LM's pre-pretraining improves the adaptation to a specific task; the evidence suggests that the generalization achieved under this paradigm can be poor because the model is overly specific to the training distribution and does not generalize well outside it [73], [74]. Additionally, training sets often do not comprehensively cover linguistic phenomena, discouraging such strategies in risk detection. Deep neural networks suffer from inductive bias overfitting to shallow syntactical features [75] or words that, as for NLI example, indicate direct contradiction [76]. Thus, the performance of fine-tuned models on specific tasks and benchmarks may exaggerate actual performance, even when it is nominally at the human level. Therefore, an effective solution shall rely on language models' generalization capabilities rather than a dedicated detection scheme. To design a risk-comprehensive system, the example-based approach shall be reduced, and instead, general principles shall be identified and explored.

Chapter 4

Proposed Solution

This chapter describes a solution to a research problem: identifying risk interaction in the narrative. In its essence, it is a document-level relationship extraction. It uses the triplet Threat-Vulnerability-Assets to normalize the risk propagation in the narrative. The propagation is represented in the Asset-Vulnerability-Hazard graph (the A-V-H graph). The Knowledge Acquisition Pipeline to construct the graph is presented. Each processing step is provided with algorithmic complexity to prove the solution's efficacy, which is multinomial as opposed to the combinatorial complexity of the naive approach. Relationship acceptance is formulated as a multicriterial optimization task. The optimization combines an ensemble of deep textual entailment classifiers and a large language model. The chapter concludes by discussing an approach to various relationship templates to construct a richer semantic representation of risk propagation.

4.1 Asset-Vunerability-Threat Triplet and A-V-H Graph

To start with the risk network model, it is required to define its main building block. The propagation of hazard can be represented by an AVH triplet (Fig. 4.1) [77]. In this representation, the *Threat* or *Hazard* impact the *Asset's* due to its interaction with the asset's *Vulnerability*.

In the triplet, the *Asset* is a component of the system, or its part, relevant from the perspective of the analysis context. *Assets* may be at different levels of abstraction. For instance, an asset may be a car at risk of an accident because of a slippery road surface and excessive speed. *Assets* will also be elements of the car, e.g., a tire prone to being punctured or engine components subject to specific failure, e.g., low oil level The *Vulnerability* is a system's component, its part, or anything that interacts with the system that impacts its performance. It causes the *Asset* to lose its ability to function correctly under the influence of a *Hazard* or *Threat*. In the case of a car, a *Vulnerability* may be a "slippery road surface" that generates the risk of an accident under the influence of a "speed" *Hazard*. In the case of an engine, a *Vulnerability* may be a "low oil level" generated by the "oil leak" *Hazard*.



Figure 4.1: Threat - Vulnerability - Asset Triple [77]



Figure 4.2: The A-V-H Graph, network approach to risk analysis. green : *Assets*, red : *Hazards*, grey: *Vulnerabilities* [1]

The "Low oil level" is contextual as it will be a *Hazard* for other *Vulnerabilities*, e.g., "high temperature."

To simplify the problem, we can assume that *Assets* are represented by nouns: "car", "pump", "engine", etc. Similarily *Vulnerabilitilies* and *Risks* are represented by nouns as well: "speed", "rain", "pressure", "temperature" etc.

The Asset - Vulnerability - Hazard, the A-V-H graph is a Knowledge Graph that is a comprehensive representation of the AVH triplets in the system domain (Fig. 4.2), [1].

4.2 Problem Statement

Given the system's description contained in the document of *D* pages, consisting of a total of *S* sentences and containing a total of *N* distinct nouns, create a three-element ordered set of nouns denoting *Assets, Vulnerability*, and *Threat* (Fig. 4.3) forming Asset-



Figure 4.3: Asset-Vulnerability-Threat Extraction Problem

Vulnerability-Threat triple (Fig. 4.1) and aggregate them into an A-V-H graph (Fig. 4.2) for a network analysis of risk interaction in the system.

In other words, *Assets* are system elements identified by *nouns* in the syntactic decomposition of the description. The *Vulnerabilities*, and *Threats* are contextual risk-related items detected for each *Asset*. They are identified by *nouns* in the syntactic decomposition as well [1], [77]. Therefore, the task is to identify the ordered triple of nouns to which we can assign and confirm the semantics of *Asset, Vulnerability*, and *Threat*. The identified triplets must be supported by the narrative from which they are extracted.

Computational Complexity of Triplet Extraction

In the document holding the description, there are N distinct nouns to be assigned to either of the categories. Estimating the upper bound of the time complexity of the task in the function of N nouns is straightforward. The function f of the number of triples in the set of N nouns is the number of combinations without repetitions of a subset of 3 elements from the set of N elements:

$$f(N) = 3! * \binom{N}{3}$$

Therefore, the problem of selecting the semantic triple is of polynomial $O(n^3)$ complexity.

The narrative must support the classification of the elements into either of the categories. However, the Entity Recognition will not allow for the direct classification of the abovementioned classes as *Vulnerability* and *Threat* are contextualized. To confirm that the narrative supports the triple, it is required to verify that there is a risk-related semantic relation among the elements in the triple. Therefore, we will select elements of the triple



Figure 4.4: Document-Level Validation Problem

et

eh

The failure to plan and account for extreme floodwaters resulted in the immediate death of 26 000 as a result of the water itself.



(Fig. 4.1) that are connected by a risk-specific relation across the document (Fig. 4.4).

Risk-Related Semantic Relation

In general, the semantic relation is defined as a triple containing (Fig. 4.5):

- *head entity* e_h which is the subject of the relation e.g. "floodwater",
- *tail entity e_t* which is the object of the relation e.g. "death",
- *predicate* which defines the kind of relation between the head and tail entities, e.g., *"has an impact on"*

The relation is detected within the textual context, which contains head and tail entities and supports the predicate between them. The risk-related relation must have properties that:

- *transitive*, meaning that if a noun *a* is a head entity e_h in relation *r* with noun *b* being a tail entity e_t and noun *b* is a head relation e_h in relation with noun *c* is a tail entity e_t , then noun *a* is a head entity e_h in relation with noun *c*. From a risk analysis perspective, it means that the hazard propagates across the system, and we can identify the impact of *Threat* through a *Vulnerability* on an *Asset*
- *irreflexive*, meaning that fact that any Asset or Vulnerability or Hazard cannot be

in relation with itself. From a risk analysis perspective, any system element cannot impact itself, and external factors must exist to initiate risk propagation. Therefore, i.e., an *Asset* cannot be a simultaneously *Vulnerability* and *Threat* to itself.

• *antisymmetric* meaning that given the context c, *a* is in the relation with *b* and *b* is in the relation with *a* only if *a* and *b* are the same entity. In other words, given the current elements in the Asset-Vulnerability-Threat triple, current *Asset* can be assigned *Vulnerability* role and current *Vulnerability* an *Asset* role if and only if there is another context that supports the switch.

The above assumptions on risk-related relations constrain the network representation of risk propagation to directed graphs without self-loops.

Triplet Validation Problem

The entities of the triple can be distributed across the entire narrative (Fig. 4.4). The AVH triple defines the risk relation as a directed relation following from *Hazard* through *Vulnerability* to *Asset* (Fig. 4.1). In the brute-force approach, the relations are checked pairwise within the fixed contexts, e.g., sentence. Therefore, we must find the combinations of contexts where entities form a risk-related relation such that a chain of relations links the triple elements.

We assume that there are N distinct nouns distributed across C contexts in the document, and there are k nouns on average in each of the contexts. The function indicating the number of comparisons to perform to confirm the existence of risk-related relation between any pair of nouns is a function of the number of contexts:

$$f(C) = (k-1)^{C-1} k^C C!$$

Therefore, the time complexity of a brute force validation is $O(k^C C!)$, and it is intractable.

4.3 Proposed Solution Architecture

There are three main issues impacting the modeling of risk propagation using the AVH triple:

- computational complexity to validate the naive approach,
- lack of training sets to use other NLP approaches, for example, apply document-level relationship extraction directly,
- no possibility to detect the AVH triple directly in the narrative using Entity Recognition approach



Figure 4.6: Proposed Solution Architecture

The solution will address these problems by:

- decomposition of the processing into specialized steps, namely, constructing Semantic Frames Graph and Intermediate Relationship Graph (IRG), which will reduce the computational complexity
- ensemble learning that will establish the acceptance threshold for the accepted relations to counter the lack of training examples and standard ROC analysis
- risk-related relation analysis will replace the entity recognition to identify AVH triple elements in the narrative
- path analysis on IRG will provide the required validation method

Overall, the solution will answer the research questions:

- First, as it is financially prohibitive, the research shall answer if it is possible to construct a risk detection system without creating a dedicated LLM for Risk Analysis.
- Second, it shall evaluate available trained and language-specialized classifiers and construct the pipeline to identify risk relations without a dedicated training set.
- Third, a validation method shall be provided.
- Finally, the pipeline shall produce the A-V-H graph, a chosen network model of risk interaction.

4.4 Sentence Decomposition and Semantic Role Labelling

Semantic Frames are semantic building blocks of sentences. They are at the heart of the *compositional semantics* that postulates that the sentence's overall meaning is a composition of its frames [78]. Semantic Role Labeling (SRL) detects frames and their structure in the sentence (Fig. 4.7). The central element identified in each frame's structure is verbs [69], for example, *offering* on (Fig. 4.7). The other elements are generalized semantic roles (ARG0, ARG1, ARGM-LOC, ARGM-TMP, etc.) [67] identified by the SRL and associated with



the frame's central verb (Fig. 4.8) We seldom encounter a situation in which an event is de-



Figure 4.7: Semantic Role Decomposition of a sample sentence

Figure 4.8: Semantic Frames Graph representation of a sample sentence

scribed in a singular sentence with a subject and object explicitly stated without additional predicates. A singular sentence would transform into a single frame, and the relationship would be expressed **directly**. However, this is not the case for complex and subordinate sentences. Therefore, when complex sentences are decomposed, then a hierarchy of semantic frames is created to capture the interaction between frames and each frame's argument separately from other frames' arguments. (Fig. 4.9).

Semantic frame decomposition is essential for relationship detection as it separates sentence components according to the frame to which they belong. This simplifies the analysis as, for example, the relation of interest may be chained in the hierarchy of frames. In a sentence *"A water landing of a jetliner that lost both engines due to hitting birds became known as the Miracle on the Hudson River"* (Figure 4.9), a subject *jetliner* is linked with *engines* and *birds* through frames *lost* and *hitting*. The main frame *known* is skipped as irrelevant from a risk analysis perspective. This simplifies the analysis of the relation between *jetliner* and other nouns - potential hazards of a jetliner:

- 1. [jetliner] "a jetliner lost both engines due to hitting birds" [engine]
- 2. [jetliner] "a jetliner lost both engines due to hitting birds" [bird]

The relation between head and tail entities is often distributed across the description. In this case, the relationship is expressed by a sequence of frames that connects them and forms a reasoning scheme that justifies the existence of the relation (Fig. 4.10). For example,





Figure 4.10: Subject - Frame(s) - Object(s)

it is possible, given the exemplary text (Fig. 3.4.), to verify that "droplet" is in a relationship with the "engine". In this case, the analysis may follow the deduction path: "droplet" \rightarrow "supercooled water droplets collide with a surface" \rightarrow "if supercooled water droplets collide with a surface they may result in blocked fuel inlet pipes" \rightarrow "blocked fuel inlet pipes" \rightarrow "aviation fuel designed for use in aircraft powered by gas turbine engines" \rightarrow "engine"

4.5 Semantic Frames Graph

Formally, a Semantic Frame Graph is an undirected, attributed, heterogeneous graph

$$G = (V, E)$$

where:

- *V* is a set of nodes of the following types:
 - Noun: nouns detected in the frame's argument. White rectangle
 - Argument: a span of text describing the semantic role of a frame. Yellow rectangle.
 - Verb: verb identifying an event associated with the frame. Green rectangle.
- *E* is a set of edges representing:

- *role_type* : a specific semantic role type for example: ARG0, ARG1, etc.
- *verb*: a connection to the central verb of the role's constituent frame. This happens if the description contains subordinate clauses.
- *noun*: a connection to a noun that is a part of the frame only.

The graph is constructed by applying recursive semantic decomposition (Algorithm 1) on every sentence in the corpus. At the sentence level, the decomposition is the recursive application of deep SRL identification [72] until none of the identified frame's elements can be further split into frames (Algorithm 1). Nouns are assigned to the lowest-level frame's role and are not repeated at the higher levels.

The results are recorded in a graph structure as a hierarchy of frames (Fig. 4.11). Connections between frames are established through nouns that frames share.

Algorithm 1 Recursive SRL decomposition	
procedure PARSE_FRAME(<i>Frame f</i> , <i>Grap</i>)	h g)
$args \leftarrow getSRL(f)$	⊳ get arguments of the frame f
$verb \leftarrow getVerb(f)$	⊳ get the verb for the frame f
$g \leftarrow add_node(verb)$	
for argument in args do	
if argument is a frame then	
parse_frame(argument,g) ⊳ furthe	er decompose current argument in case it is a
frame	
else	
$g \leftarrow add_link(verb, argument)$	▷ connect verb's frame with its argument
for n in getNouns(argument) do	•
$g \leftarrow add_link(argument, n)$	▷ connect frame's argument with its nouns
end for	
end if	
end for	
end procedure	


Figure 4.11: Semantic Frame Graph Structure

Computational Complexity

Given that, on average, a sentence has f frames, the decomposition of S sentences into its frame structure is a function of a number of sentences f(S) = f * S, which is linear with the number of sentences.

4.6 Relationship Extraction and Semantic Pattern

The semantic frame defines a deep case as an atomic semantic element [78]. It defines the actor (subject), the object, and the elements. A single frame defines the semantic relationship between its arguments.

Within the document, the relations between two adjacent nouns are detected by applying a semantic pattern on the path of frames connecting them in the SFG. The semantic pattern, a template *t* for the relationship context *c* is a textual core of the relation that requires providing *subject* i.e., *head entity* (e_h) and *object* i.e. *tail entity* (e_t). The substitution is called *verbalization* (Fig. 4.12) [48].

The frame defines the relation context, which should support the verbalization. Entailment evaluates how well it is supported. The premise is the relation context, and the



Figure 4.12: Semantic Pattern Verbalization: template t, relation context c

hypothesis is the template verbalization.

For example, the frame: "The failure to plan and account for extreme floodwaters resulted in the immediate death of 26000 as a result of the water itself" contains one verb, "resulted" suggesting some form of impact (Fig. 4.12). It is a premise. Assuming that the *t*: 'has impact on', head entity: e_h : "floodwater" and tail entity: e_t : "death", then the verbalization of the hypothesis is "floodwater has an impact on death". Therefore, the relation detection is cast as NLI problem [47].

Relationship extraction over the SFG graph verbalizes a relationship template *t* on frames connecting every *Noun* node in the graph. The structure of the SFG graph allows traversing it using an explicitly structured walk on the graph, namely, metapaths [79] with a limited number of metapaths:

- 1. noun-role-verb-role-noun where the relationship is contained within a single frame
- 2. *noun-role-noun* where the relationship is contained within the frame's argument, which is a special case of a single frame,
- 3. a: *noun-role-verb-role-noun-role-noun*, b: *noun-role-verb-role-noun-role-verb-role-noun*, where the relationship is held in adjacent frames
- 4. *noun-role-verb-role-verb-role-noun* where the relationship is contained in two frames of a subordinate clause

.The structured walks allow the detection of relationships specific to each metapath. There are two cases (Fig. 4.13):

- single frame: metapaths 1 and 2,
- two frames: metapaths 3a, 3b and 4.



Figure 4.13: Relation Detection per Metapaths



Figure 4.14: Single Frame Case



Figure 4.15: Two Frames Case

4.6.1 Single Frame Relation Extraction

A single frame case is a basic case in which the relationship between the head and tail entities is detected within the frame. This is template verbalization case only [48], e.g., *water vapour create droplets* is validated against a relationship template *T*: "has effect on", e.g., *"water has impact on vapour"* (Fig. 4.14).

Computational Complexity

As sentences have, on average, k nouns, in the worst-case scenario, all k nouns would be part of each frame. We have S sentences in the document, which contain, on average, f frames; therefore, the number of comparisons to evaluate template entailment is given by the linear function of sentences in the document: f(S) = k * (k - 1) * f * S. The time complexity is then linear for the number of sentences in the description: O(S)

4.6.2 Two-frame Relation Extraction

Written sources of risk propagation often assume a correlation of entailment with other entities based on discourse coherence. For example, in the following sentences *"Electricity was cut off in the control room. All electronic equipment in the room was disabled."* it is natural for a human to infer that *"Electricity has an effect on electronic equipment in the control room"*. A human reader naturally considers the example as coherent as the *discourse focus* does not fluctuate [80].

The two frames scenario targets the case directly. Additionally, the scenario will combine frames that may not follow the passage's "reading" order, which allows general crossdocument relationship detection. However, having a single noun in common does not make the frames linguistically coherent. It takes more than a common noun from the case above, i.e., *"room"*. For frames to be linguistically coherent, they must form a semantic flow, namely, they must depict a discourse or dialog coherence [80], [81].

The driving assumption behind the centering approach is that "focus", the most salient entities discussed, smoothly *transition* across the discourse. It uses the notion of *center* to measure the focus shift between two consecutive sentences. The change is measured between the "forward-looking centers", C_f , which are nouns in the preceding sentence, and "backward-looking centers" C_b [80] which are nouns in syntactically prominent (subject or object) positions in the succeeding. The centering theory defines three categories of focus transition: "*continuation*" where the preceding subject is repeated, "*retaining*" where the grammatical role between the elements in C_f and C_b changes and "*shift*" where elements in the centers are different.

The centering theory was developed in the 90s; therefore, one of its limitations is that it can only consider surface forms of centers and does not account for synonyms or contextual representations of words. Consequently, we retain the original assumption on prominent roles of centers (subject and objects); however, we measure the *"semantic drift"* which is how the semantics of the succeeding sentence C_b is represented in the preceding dropping the analysis of forward-looking centers overlap entirely.

The verbalizer approach cannot be used directly on both frames nor the complete path as it does not comply with NLI training [47]. The relationship between entities (nouns) exists if a path in SFG joins them [40]. A path of semantic frames joining them creates a passage defining a relationship. For example, a path between *water* and an *engine*: "*water*" \rightarrow "supercooled water droplets collide with a surface" \rightarrow "if supercooled water droplets collide with a surface they may result in blocked fuel inlet pipes" \rightarrow "blocked fuel inlet pipes" \rightarrow "aviation fuel designed for use in aircraft powered by gas turbine engines" \rightarrow "engine"

The generated path may indicate that *"water"* impacts *"engine"*. However, we cannot 'entail' the whole path as it would not comply with the training scheme for entailment, which considers inference on a single sentence only [47] [65]. Therefore, the relationship detection task is split into coherency and single-frame template entailment and performed pair-wise for each consecutive frame in the path. The metapath walks 3a, 3b, and 4 perform the required pair-wise validation for each node's adjacent frames in the SFG graph.

4.6.3 Modified Dialog Coherence Function

Let F_i and F_{i+1} be the preceding and succeeding frames in the path. The backwardlooking centers of the pair of frames F_i and F_{i+1} are generalized a subject or an object of the frame F_{i+1} Generalized subject and object nouns are linked to ARG0, ARG1, ARG2, ARG4 roles in the SFG decomposition of the frame F_{i+1}

Let $f_{i,i+1}^c$ denote a coherence function between both frames F_i and F_{i+1} .

The coherence functions shall be bounded 0 $f_{i,i+1}^c$ 1 such that the interpretation if, from the linguistic standpoint, is that if $f_{i,i+1}^c \approx 0$ means that both frames are not coherent and if $f_{i,i+1}^c \approx 1$ are coherent.

The function shall evaluate the following pair of frames (Example 1): F_1 : *water vapour create droplet* and F_2 : *spark causes ignition of fuel* as incoherent.

The following pair (Example 2): and F_1 : water vapour create droplet and F_2 : droplet

can block the fuel inlet pipe as more coherent than Example 1.

To calculate the coherence, we will use zero-shot text classification TC [82] of the backward-looking centers C_b in the frames and normalizing the output, hence defining Modified Dialog Coherence Function to be:

$$f_{i,i+1}^{c} = \frac{1}{|C_b|} \sum_{m=1}^{|C_b|} \left\{ \frac{1: TC(F_i, m) > TC(F_{i+1}, m)}{TC(F_i, m) / TC(F_{i+1}, m): TC(F_i, m) / TC(F_{i+1}, m)} \right\}$$
(4.1)

The function (4.1) defined as such has the properties:

- $f_{i,i+1}^c \approx 0$ if there is no reference to any of the backward centers in F_i ; hence frames are incoherent. The score for Example 1 is 0.0022
- *f*^c_{*i*,*i*+1} ≈ 1 if there is a perfect overlap of the centers in both frames, meaning both frames are almost identical semantically.
- $0 < f_{i,i+1}^c < 1$ if there is an overlap of the centers and frames are semantically related. The score of Example 2 is 0.6

Computational Complexity

The metapaths constrain the way the SFG graph is traversed. The worst-case scenario requires the most computation when all relationships are represented as metapath 3b, which connects nouns separated by two frames sharing one noun (Fig. 4.16). Additionally, in the scenario, each sentence in the text is singular and decomposed to a single frame; therefore, each frame has an average number of *k* nouns. To evaluate all relationships, we must visit all metapaths 3b subgraphs and traverse them for (k - 1) * (k - 1) pairs. For each pair, we have to find the shortest path connecting them using the Dijsktra algorithm - complexity O((V + E)log(V)), where V is the number of vertices in the subgraph. The subgraph is sparse $V^2 >> E$, therefore the complexity has the form: O(Vlog(V)). The number of vertices is easy to estimate as there are 2k nouns, 2k roles, and 2 verbs. Therefore, the time complexity to evaluate all relations in a metapath 3b can be estimated as

$$O(k) = (k-1)^{2}(4k+2)\log(4k+2)$$

depends on the average number of *k* nouns per sentence only. In the worst-case scenario, all metapaths 3b will be connected. Therefore, the upper bound on complexity is the number of pairs of frames, which is proportional to the number of sentences

$$O(S) = S * (S - 1)$$



Figure 4.16: Metapath 3b: the worst-case scenario

4.6.4 Combined Entailment Function

Combined entailment $f_{i,i+1}^e(t)$ measures entailment of relationship template t within two adjacent frames in the path F_i and F_{i+1} . It evaluates the entailment of the template of the relationship in both frames and their dialog consistency. The function is bounded: 0 $f_{i,i+1}^e(t)$ 1.

Let n_i be a noun in the path leading to frame F_i like "temperature" in example (Fig. 4.15), n a noun linking F_i and F_{i+1} like "water" in example (Fig. 4.15), n_{i+1} a target noun like "vapour" in example (Fig. 4.15). Let t denote the relationship template "has an effect on", which would mean that $f_{i,i+1}^e(t)$ would measure if "temparature has an effect on vapour" in two frames setup (Fig. 4.15). Let $t(n_1, n_2)$ denote verbalization of the template given pair of nouns (n_1, n_2) , for example, t("temperature", "water) would resolve to "Temperature has an effect on vapour".

$$f_{i,i+1}^{e}(t) = min[f_{i,i+1}^{c}, max[RTE(F_{i}, t(n_{i}, n)), RTE(F_{i+1}, t(n, n_{i+1}))]]$$
(4.2)

where RTE is a classifier detecting *entailment* classification probability given frame *F* and verbalization of the template *t*. We used transformer RTE implementation [48].

Following the example, we calculate the value of the RTE of the template *t* "has an effect on" verbalized as "temperature has an effect on water" of the frame "Low temperature and high altitudes cause water vapour to create droplets"; the template verbalization "water has an effect on vapour" for the frame "water vapour create droplets"; dialog consistency between both frames and apply the rule (4.2).

4.6.5 Multiple Templates

The goal of the solution is to model the propagation of hazards. This defines the requirements for semantic relations to be:

- irreflexive
- transitive
- antisymmetric.

Only such relations will model the potential risk propagation within the system. For example, a template *has effect on* meets the criteria as no object can affect itself without external cause. The effect propagates, meaning that if component a affects another component b and the component b affects another component c, then component a affects component c. The effect relation is antisymmetric, meaning that given narrative justifying the effect between components a and b, such as *"airplane uses engines for flying"* does not automatically mean the opposite relation unless detected in the narrative specifically. The example justifies the hazard propagation from *engine* to *airplane* only, not vice versa. If the narrative is expanded with the text: *The engine's performance relies on fuel stored in the airplane's wings.*, only then the *has effect on* relation can be established between *fuel, airplane* and *wing* and *engine*.

There are more ways to express the hazard propagation. The hazard will likely propagate across the structure if an element is "a part of" another element. The "a part of" relation meets the criteria for risk-related relation as it is irreflexive - an element can be a part of itself, transitive and antisymmetric - if an element is a part of the other, the other is a *whole* not a part of the first. For example, a "battery" is a part of the "engine". Hazards associated with it will propagate onto the "engine", but not necessarily vice-versa. Another example is a "a type of" relation, which establishes the taxonomic structure of the elements in the narrative. The "a type of" relation meets requirements for risk propagation relation as well as it is irreflexive (hardly any object can be a type of itself), transitive (if an object a is a type of the object c). A "a type of" relation is antisymmetric by definition. Continuing with the fuel example, if the narrative is expanded with "The ATF is a type of aviation "fuel," then "ATF" connects "fuel" and propagates hazard to elements "fuel" would propagate: "engine" and "airplane".

From the risk propagation modeling perspective, *precise* detection and comprehensive representation of semantic relations between objects is not required. For example, another risk-related template is introduced *"is a part of"* with the context "Engine is a part of the airplane". From a risk propagation perspective, it is irrelevant if the risk propagation between "engine" and "airplane" is established through *"is a part of"* or *"has effect on"*. Therefore, selecting the best risk-related relationship between the entities the narrative supports, namely, with the highest coherence score, is enough.

Formalizing the approach, let $E_{i,j}^k$ denote the set of edges in the IRG graph connecting nouns n_i, n_j for the template k. Let T denote the set of templates. Let $w_k(i, j)$ denote a weight, dialog score, between nodes for the template k. Then the relationship $R_{i,j}$ connecting these nodes is defined as such a template t that maximizes the total score between the nouns n_i , n_j across all the templates and metapaths. Select the best-supported template:

$$R_{i,j} = \operatorname*{argmax}_{t \in T} w_t(i,j)$$

4.7 Intermediate Relationship Graph

The Intermediate Relationship Graph (IRG) is a Semantic Network that stores all detected relations between nouns by applying metapath walks on the SFG. Formally, the IRG a weighted multigraph (Fig 4.17), where nodes are nouns and edges represent the verbalizations of the template of, for example, *"has effect on"* relation; edges' weights are either entailment score or value of the Combined Entailment Function depending on the walk; the metapath that connects the nouns in the original SFG graph is recorded as an edge attribute.

Both verbalized relation and dialog consistency are transitive. Therefore, it is possible to traverse the IRG freely, meaning the NLI problem has been transformed into a graph traversal one. For example, node n_1 is directly connected with node n_3 , which indicates a direct impact. It is also connected through node n_2 , which could indicate a mediated impact and form a *a chain of impact*.

While traversing the graph, we can use only edges with scores above the assumed threshold and analyze the evidence. For example, there is a metapath joining nodes n_1 and n_2 (Fig. 4.17) in the SFG graph. Weight w_3 is the score of the template verbalization of the connection, and m_3 means, depending on the metapath, either a frame or two frames connecting them.

From the risk analysis perspective, the graph provides a complete propagation of impact between system elements, where edges provide scores and evidence.

The IRG graph is a central point of representation of risk within a single document and across documents. It allows incremental knowledge acquisition as new relations can be added to the graph.

Computational Complexity

The time complexity of constructing the IRG graph from SFG is bounded by the complexity of the most complex metapath (3b), and it is bounded by the number of sentences



Figure 4.17: Intermediate Relationship Graph Structure. Attribute "m": metapath, "w" - weight, "isYes" - LLM's validation results

in the narrative, not the number of nouns:

$$O(S) = S^2$$

4.8 Asset-Vulnerability-Hazard Graph

The Asset - Vulnerability – Hazard (A-V-H) graph aggregates the relationships expressed in the IRG graph that meet the minimum weight criterion for the weights¹. It isolates three types of nodes essential from the system analysis perspective, namely:

- Asset node: an element of the system that is relevant from the perspective of its proper operation in the analysis context. Assets may be at different levels of abstraction. For instance, an asset may be a car at risk of an accident because of a slippery road surface and excessive speed. Assets will also be elements of the car, e.g., a tire prone to being punctured. Assets will be engine components subject to specific failure, e.g., low oil level
- Vulnerability node: an element of an object, an event, or any other object that causes the Asset to lose its ability to function correctly under the influence of a Hazard. In the case of a car, a Vulnerability may be a "slippery road surface" that generates the risk of an accident under the influence of a "speed" Hazard. In the case of an engine,

¹the threshold is a solution to a multicriterial optimization task



Figure 4.18: A-V-H Graph Structure

a Vulnerability may be a "low oil level" generated by the "oil leak" Hazard. "Low oil level" will be a Hazard for a Vulnerability, e.g., "high temperature."

• Hazard node: a system element, an event, or other asset that exposes the Asset to a risk due to the Vulnerability

Relying on the transitivity of relations in the IRG graph, the construction of the specific A-V-R nodes is performed as follows:

- Asset nodes will be all IRG nodes by default.
- Vulnerability modes will IRG nodes that directly affect Assets. Hence, these are the nodes that directly neighbor with Assets.
- Hazard nodes will be IRG nodes that connect with the Asset through Vulnerability. These are nodes from which the Asset node is reachable through its Vulnerability.

Therefore, the aggregation transforms a weighted, directed multigraph IRG into an unweighted, directed graph A-V-H (Fig. 4.18). The connection between nodes in A-V-H is established if the IRG graph connects them with an edge whose weight is above the calculated threshold or if there is a path between nodes connected with all edges' weight above the assumed threshold.

The A-V-R graphs enable risk analysis using a graph analysis approach, i.e., nodes'

centrality measure, to identify hazard-related hubs that aggregate and transmit influence to other system elements [1].

Computational Complexity

The time complexity to establish *Asset - Vulnerability* pair is given by the time complexity of evaluating neighbors of *Vulnerability* nodes. In the IRG, sentences and frames are not directly represented. They are a part of the edge attribute - the metapath connecting the nouns. In the worst-case scenario, when IRG is a complete graph, each noun is connected in the narrative by single or two-frame metapaths. Each node will have N - 1 neighbors in this case. Therefore, the upper bound on the time complexity to identify the pair is

$$O(N) = N * (N - 1)$$

The relation between *Vulnerability* and *Risk* exists if a path links them in the IRG. Therefore, the time complexity is estimated by the complexity of finding the shortest path. The shortest path, however, shall skip the *Asset* node as it is already connected with the *Vulnerability*. Therefore, the complexity of a single path between selected *Asset* and *Vulnerability* is O((N-2)log(N-2)). Again, in the worst-case scenario, we must scan all IRG nodes for *Asset* and *Vulnerability* pairs to find the shortest paths to all remaining nodes. Therefore the overall time complexity is

$$O(N) = N * (N-1) * (N-2) log(N-2)$$

which is

 $O(N^3)$

4.9 Validation and Explainability

4.9.1 Validation and Threshold Calculation

Validation in the context of machine learning and data science means evaluating the model on unseen data. In the proposed solution, as no specific dataset holds unseen data, validation is performed by another model, the LLM, performing the same relation extraction task. It aims to establish the cutoff threshold for the weights on the IRG graph. The *threshold* is the value of the dialog function that cuts off a maximum number of relations rejected by



Figure 4.19: Validation Approach

LLM, keeping the maximum number of relations confirmed. The approach is based on the ensemble of, in fact, three independent classifiers performing the same Language Inference task ((Fig. 4.19):

- 1. the entailment and text classification that are combined to provide the estimation of the *coherence function*, which is a *weight* on the IRG edges,
- 2. the prompted LLMs, which provides another independent decision on whether the edge encodes the relation template and the edge is justified. LLM works on prompt engineering principle in a zero-shot setup [83], [13], [84] which is encoded as a *isYes* parameter on the edge. The value of isYes: 1 indicates that LLM confirms the verbalization, and 0 indicates otherwise.

The LLM prompt is constructed to provide explicit instructions on the task: "Given the premise text only decide if given the premise, the hypothesis is true. Answer Yes if true or No if unclear or false.". Answer 'Yes' is converted to 1, 'No' to 0 and recorded as an edge attribute (Fig. 4.17).

Formulation of the Optimization Task

Let $S = \{s_1, s_2, ..., s_n\}$ be the set of elements denoting distinct *weights* in the IRG graph with additional attributes associated with them and each element s_i in the set *S* contains:

- *w_i*: the weight value,
- $y_i^{(0)}$: number of edges having weight w_i in the IRG classification unconfirmed by LLM,
- $y_i^{(1)}$: number of edges having weight w_i in the IRG classification confirmed.

The elements of the set *S* have the property that they are indexed according to the *weight*:

$$\forall s_i, s_j \in S$$
, if $i > j$ then $w_i > w_j$

Let

- $t^{(0)}$ denotes the total number of unconfirmed relations in the IRG and
- $t^{(1)}$ denotes the total number of confirmed relations in the IRG.

Then, we can define functions:

• $f^{(0)}(w)$ returning a fraction of cumulative unconfirmed relations in IRG below or at w

• $f^{(1)}(w)$ returning a fraction of cumulative confirmed relations in IRG below or at wThe values of the functions are calculated as:

$$f^{(0)}(w) = \frac{\sum_{i=1}^{n} y_i^{(0)} | w_i <= w}{t^{(0)}}$$
$$f^{(1)}(w) = \frac{\sum_{i=1}^{n} y_i^{(1)} | w_i <= w}{t^{(1)}}$$

Let

$$\mathbf{F}(w) = (-f^{(0)}(w), f^{(1)}(w))$$

Then, we can formulate the multi-objective optimization task to evaluate the IRG edge cutoff threshold of w_{th} , below which we will maximize the number of unconfirmed and minimize the number of confirmed relations to be removed from the IRG.

$$\min_{w} F(w) \tag{4.3}$$

The optimization task Eq. 4.3, can be solved by selecting the element of the Pareto Front: $(f^{(0)}(w_{th}), f^{(1)}(w_{th}))$ such that

$$\min_{w} \|\mathbf{F}(w) - \mathbf{z}^{ideal}\|$$

where

$$\mathbf{z}^{ideal} = (0,0)$$

4.9.2 Explainablity

Explainability in AI refers to the possibility that a human can understand, interpret, and accept decisions made by artificial intelligence agents. It goes past performance scores, e.g., ROC curve [85] or measures such as dialog coherence or textual entailment. In the proposed solution, the IRG edges store the metapath, dialog coherence score, and relation template used to detect the relation between the designated SFG nodes. The explainability



Figure 4.20: Explaiablity Chain. The path of nouns connected by the IRG edges with their attributes.

'Path': [('liver', 'noun'),('clotting', 'noun'),('fibrinogen', 'noun')]
Edge: ('liver', 'noun') - ('clotting', 'noun') Properties: {'frames:': 'the liver synthesizes clotting factors such as fibrinogen prothrombin.', 'weight' : '0.9378', 'label': 'has effect on'}
Edge: ('clotting', 'noun') - ('fibrinogen', 'noun') Properties: {'frames:': 'the liver synthesizes clotting factors such as fibrinogen prothrombin.', 'weight' : '0.8299', 'label': 'is a type of'}

Figure 4.21: Explainability Example

means that the solution can provide the exact chain of frames connecting the required nodes (Fig 4.20), providing all required information supporting the existence of the specific relation between them. For example, the sample narrative describing the function of the liver is analyzed to detect how the liver impacts other human body mechanisms, such as hemostasis. Such an approach would allow us to evaluate the potential impact of the medicine, which mechanism of action targets the liver. This functionality is important as medicine can indirectly affect other essential processes. The example (Fig. 4.21) explains how the liver impacts fibrinogen, given the narrative on the liver.

4.10 Summary

There has been limited success in designing a Knowledge Acquisition Pipeline for a comprehensive network representation of risk interactions. Unfortunately, most current solutions rely on the dedicated detection schema targeting either a specific input format [31] or a specific domain [6]. The proposed solution relaxes the training set constraints. It applies existing and available deep NLI classifier BART [82] to infer the relation through the entailment between relationship context expressed as the combination of frames.

Compared with the direct approach to identifying risk-interaction triples (Asset-Vulnerability-Risk) in the narrative, the proposed detection and validation are performed in multinomial time complexity, which is feasible for large descriptions. The proposed solution achieves all the research goals given:

Large Language Models

Although it is possible to train a dedicated LLM, i.e., LLama [86], several challenges

would have to be resolved to use such LLM directly:

- insufficient textual data
- still expensive infrastructure
- a prompting strategy to identify risk-related relations.

Current general LLMs, such as chatGPT, LLama, or FLAN, encode general language knowledge. Their general linguistic capabilities were used, i.e., to confirm effect propagation instead of threat or risk directly.

• Existing general-purpose language classifiers

Alternatively, we can train a dedicated transformer-based [35] classifier to classify all risk-related candidate relations in the narrative. Such an approach, however, would require a significant training set, which is not available for the risk domain. Instead, the solution suggests using a verbalized semantic template entailment approach that scores (values of the Dialog Function) are the plausibility of the template expressed in the semantic frame.

• Validation Approach

The relationship detection results are validated twofold:

- through multiobjective optimization, which establishes the acceptance threshold for entitlement scores, - through a direct "reasoning" path: the sequence of frames with the Dialog Function scores that connect the entities.

• Knowledge Acquisition Pipeline for A-V-H graph

The processing constructs the required network representation of risk interaction, namely the A-V-H graph.

Chapter 5

Results

In this chapter, I discuss the efficacy of the proposed solution, aiming to provide an analysis of its performance and impact. Building on previous chapters' theoretical foundations and design principles, I focus on empirical evaluation and practical implementation. This chapter is a critical juncture in my work, transitioning from conceptual frameworks to practical results.

The solution's efficacy will be discussed using the DocRED dataset, which is used as a benchmark for general-purpose inter-sentence detection. The multi-template example will show that it is possible to construct more comprehensive IRG and A-V-H graphs. Lastly, the impact and selection of various language models will be discussed.

In addition, synthetic examples will explain the functionality in greater detail. I also explore qualitative outcomes through case studies and practical applications. These realworld examples illustrate the solution's versatility and potential to address various challenges across domains.

This chapter aims to substantiate the proposed solution's value, demonstrate its capacity to deliver meaningful improvements, and lay the groundwork for future advancements.

5.1 Knowledge Acquisition Pipeline

The synthetic example of the risk-related report contains a description of the behavior of the fuel and its impact on the engine and aircraft (Fig. 5.1). Although it is not explicitly expressed, there is a logical connection indicating that potential risk interaction between "water" and "fuel" in the context of the "airplane" exists. The SFG decomposition of the description (Fig. 5.2) confirms that such a path exists:

"water" \rightarrow "supercooled water droplets collide with a surface" \rightarrow "if supercooled water droplets collide with a surface they may result in blocked fuel inlet pipes" \rightarrow "blocked fuel inlet pipes" \rightarrow "aviation fuel designed for use in aircraft powered by gas turbine engines" \rightarrow "engine".

It also shows that the path is relevant from a risk interaction perspective, i.e., "water"



ATF is a type of aviation fuel designed for use in aircraft powered gas-turbine engines. If these supercooled droplets collide with a surface they can freeze and may result in blocked fuel inlet pipes."



Figure 5.2: ATF fuel SFG decomposition

- "collide" - "with a surface," which means that both water and surface affect each other. There is a similar relationship between "droplets" and "surface". In the second element, the relationship between "droplets" - "results" - "pipes" is also a "effect type" relationship between "droplets" and "pipe". The "blocked fuel inlet pipes" are connected to the last frame, where it is possible to identify a relationship between "engine" - "use" - "fuel". By reading the entire path, it is also possible to confirm that, in terms of consistency and general operational safety, the "droplet" (Hazard) affects the "engine" (Asset) through "fuel" (Vulnerability).

e The template used to detect the risk flow in the system will use the "has an effect on" pattern. Such a relation meets the requirement for risk-related relations as it is irreflexive (as the system can hardly affect itself) and transitive as "effect flows" through the system. The effect flow is antisymmetric, as the direction of the flow must be supported explicitly



Figure 5.3: Distribution of relations per detection function assigned to a metapath



Figure 5.4: Distribution of valid relations per detection function assigned to a metapath

by the relation context.

All specified metapaths perform the SFG graph (fig. 5.3):

- isRTE: indicate that metapath 1 and metapath 2 and these are intra-frame cases,
- isDialogRTE: indicate metapth 3a and 3b, which are inter-frame cases,
- isDialogRTE2: indicate metapaht4, a subordinate clause decomposition case.

The relation verbalization statistics shows that most of the relations evaluated were *isRTE* (over 50%), meaning they were identified within a single frame. However, dialog-based strategies denoted as *isDialogRTE* for two adjacent frames and *isDialogRTE2* for subordinate clause decomposition, together add a significant part of them (Fig. 5.3). In many cases, they are irrelevant as their dialog coherence score is low (Fig. 5.5), but not considering them at all would deteriorate the risk model as they are a significant part of confirmed relations (Fig. 5.4). A multicriterial optimization performed for chatGPT 3.5 LLM (Eq. 4.3) establishes the threshold $w_{th} = .18$, and relations accepted will have dialog

source	target	function	GPT Decision	weight
('airplane', 'noun')	('aviation', 'noun')	isDialogRTE	Yes	0.5047957187956318
('airplane', 'noun')	('fuel', 'noun')	isDialogRTE	Yes	0.5047957187956318
('airplane', 'noun')	('type', 'noun')	isDialogRTE	No	0.5688557397630128
('airplane', 'noun')	('use', 'noun')	isDialogRTE	No	0.0
('airplane', 'noun')	('gas', 'noun')	isDialogRTE	No	0.0
('airplane', 'noun')	('aircraft', 'noun')	isDialogRTE	Yes	0.0
('airplane', 'noun')	('turbine', 'noun')	isDialogRTE	No	0.0
('airplane', 'noun')	('engine', 'noun')	isRTE	No	0.7587770819664001

Figure 5.5: Sample of relations detected for the 'airplane'

source	target	path	verbalizer	score	prompt results
droplet	pipe	these supercooled droplets collide with a surface. If these supercooled droplets collide with a surface they may result in blocked fuel inlet pipes.	droplet has an effect on pipe	0.7	No
airplane	fuel	An airplane uses engines for flying.aviation fuel designed for use in aircraft powered gas turbine engines.	airplane has effect on fuel	0.5	No
fuel	airplane	An airplane uses engines for flying.aviation fuel designed for use in aircraft powered gas turbine engines.	fuel has effect on airplane	0.9	Yes
fuel	pipe	If these supercooled droplets collide with a surface they may result in blocked fuel inlet pipes	fuel has effect on pipe	0.6	Yes
pipe	fuel	If these supercooled droplets collide with a surface they may result in blocked fuel inlet pipes	pipe has effect on fuel	0.6	Yes

Figure 5.6: Examples of classification decisions

scores above it (Fig. 5.7). It is worth noting that the prompted chatGPT 3.5 itself does not provide perfect decisions, and solely relying on the LLM's Yes/No answer will not produce a trustworthy risk model.

Incremental Knowledge Acquisition Pipeline

The synthetic example can be expanded with additional contextual information. This simulates a scenario when the risk model is created gradually once new descriptions are



Figure 5.7: The distributions of chatGPT prompted decision together with the threshold (red vertical line) (left). Empirical cumulative distribution of relations wrt dialog function score (right)



Figure 5.8: The IRG for the synthetic example



Figure 5.9: The A-V-H graph for the synthetic example.Blue nodes: Assets, Orange: Vulnerabilities, Red: Risks

available. A new SFG graph is constructed for a new description (Fig. 5.11). The IRG graph is *augmented* with incoming new relations. Eventually, the A-V-H graph will be updated with new interactions.







Figure 5.11: SFG decomposition of additional synthetic context

The augmented A-V-H graph shows that *fuel* Vulnerability becomes a hub that connects with eight risks in the context of the *airplane* Asset.



Figure 5.12: Augmented IRG Graph for the synthetic example. New connections are highlighted.



Figure 5.13: Augmented A-V-H Graph indicating new Risks: "temperature" and "altitude"

Adding additional contextual information enriches the A-H-V graph further. For example, adding elements directly impacted by fuel, like the APU, air conditioning, electricity, navigation instruments, and others, is possible. The centrality analysis of Vulnerability nodes can help identify areas requiring specific precautions given connected risks [1].

5.2 Intra-Sentence Relation Detection

Although no available resources could be used to evaluate the proposed detection quality in document-level scenarios in the risk domain, at least one targets a similar scenario for a general case. DocRED is a large human-annotated dataset based on Wikipedia and WikiData. It is annotated for entity types and relations between them [49]. The entity types are essential factors in improving relationship detection. For example, a relation such as *"is located in"* would be possible between *CORPORATION* and *LOCATION* entities rather than *PERSON* and *LOCATION* [57]. The annotation provides examples of the relation between the entities and the sentences supporting it. Each test case is provided with its specific multi-sentence passage, for example, Zest Airplines description for (Fig. 5.14) and annotation (Fig. 5.15).

Zest Airways, Inc. operated as AirAsia Zest (formerly Asian Spirit and Zest Air), was a low - cost airline based at the Ninoy Aquino International Airport in Pasay City, Metro Manila in the Philippines . It operated scheduled domestic and international tourist services, mainly feeder services linking Manila and Cebu with 24 domestic destinations in support of the trunk route operations of other airlines. In 2013, the airline became an affiliate of Philippines AirAsia operating their brand separately. Its main base was Ninoy Aquino International Airport, Manila. The airline was founded as Asian Spirit, the first airline in the Philippines to be run as a cooperative. On August 16, 2013, the Civil Aviation Authority of the Philippines (CAAP), the regulating body of the Government of the Republic of the Philippines for civil aviation, suspended Zest Air flights until further notice because of safety issues . Less than a year after AirAsia and Zest Air 's strategic alliance , the airline has been rebranded as AirAsia Zest . The airline was merged into AirAsia Philippines in January 2016.



Zest Airways, Inc. : headquarters location : Pasay City. Evidence sentence(s): 0 Zest Airways, Inc. : country : Philippines. Evidence sentence(s): 2,4,7 Zest Air : country : Philippines. Evidence sentence(s): 6,7 Pasay City : country : Philippines. Evidence sentence(s): 0 Pasay City : located in the administrative territorial entity : Metro Manila. Evidence sentence(s): 0
Philippines : contains administrative territorial entity : Metro Manila. Evidence sentence(s): 0
Manila : country : Philippines. Evidence sentence(s): 0,3
Metro Manila : contains administrative territorial entity : Pasay City. Evidence sentence(s): 0
Metro Manila : located in the administrative territorial entity : Philippines. Evidence sentence(s): 0,3
Metro Manila : country : Philippines. Evidence sentence(s): 0,3
Ninoy Aquino International Airport : located in the administrative territorial entity :
Pasay City. Evidence sentence(s): 0,3
Ninoy Aquino International Airport : country : Philippines. Evidence sentence(s): 0,3 Asian Spirit : country : Philippines. Evidence sentence(s): 4

Figure 5.15: DocRED multi-sentence annotation example: Zest Airlines. [49]

DocRED specifies 96 relations, but not all meet the risk-specific irreflexive, antisymmetric, and transitive criteria, for example, relation: 'P54' - 'is member of sports team'. However, DocRED contains some that meet them, such as location-specific relations. They are irreflexive, transitive (e.g., Ninoy Airport is located in Manila, Manila is located in the Philippines, then Ninoy Airport is located in the Philippines), and antisymmetric (Manila is located in the Philippines, but the Philippines is not located in Manila). Therefore, they can be used as a proxy for risk-type relations. The satisfactory performance of the proposed solution in location-specific relation detection will indicate the general quality of risk-related scenarios.

The results achieved for a sample (Fig. 5.14 are presented. The location hierarchy, with

the Philippines being the top level with no outgoing edges (Fig. 5.16), suggests that the chain of location relations is preserved. Although compound nouns such as "Ninoy Aquino International Airport" are split into individual nouns: "ninoy", "aquino" and "airport" each of them is connected with its location individually e.g. "passay" - "manila" and "passay" - "philippine" as indicated in the training set. However, "manila" relations (Fig 5.17), suggest that there are additional elements detected.

Detection performance is evaluated as follows. The number of possible relations N = 1406 is the number of pairs of nouns: head and tail entities connected in the SFG graph and, therefore, have a path between them in the IRG. The total number of detected relations (106) are those whose score is above the calculated threshold. The total number of true positive (*TP*) relations is the number of pairs confirmed by the training set. The number of false positive (*FP*) relations is the difference between the number of relations detected and the true positive (*TP*) number. The number of false negative (*FN*) relations is the number of true number of true negative (*TN*) relations is the difference threshold. The number of true negative (*TN*) relations is the difference between threshold. The number of true negative (*TN*) relations is the difference between total *N* and true positive (*TP*).

Detection performance is calculated through a confusion matrix (Tab. 5.1) and standard metrics are as follows:

• Accuracy:

$$\frac{TP+TN}{TP+TN+FP+FN} \approx 0.94$$

• Precision:

$$\frac{TP}{TP + FP} \approx 0.21$$

• Recall:

$$\frac{TP}{TP + FN} = 1$$

• Specificity:

$$\frac{TN}{TN+FP} \approx 0.98$$

• F1:

 $\frac{2*Precision*Recall}{Precision+Recall} \approx 0.34$



Figure 5.16: Zest Airways "is located at" results: Philippines (philippine) focus



Figure 5.17: Zest Airways "is located at" results: Manila focus

		Pred	icted
		True	False
Actual	Positive	23	83
Actual	Negative	1383	0

Table 5.1: Zest Airways, "is located at" confusion matrix

relation	detection results	correct?
Zest Airways, Inc. : headquarters location : Pasay City.	Passay and City reachable from node 'airline'	yes
Zest Airways, Inc. : country : Philippines.	Philippines is reachable from nodes: zest, air, airline	yes
Zest Air : country : Philippines.	Philippines is reachable from nodes: zest, air, airline	yes
Pasay City : country : Philippines.	philipine' node is a neighbor of 'passay'	yes
Pasay City : located in the administrative territorial entity : Metro Manila.	manila' and 'metro' nodes are neighbors of 'passay'	yes
Philippines : contains administrative territorial entity : Metro Manila.	This is the opposite to 'is located in' relation. It is not covered by the template	no
Manila : country : Philippines.	philipine' nodes is a neighbor of 'manila'	Yes
Metro Manila : contains administrative territorial entity : Pasay City.	This is the opposite to 'is located in' relation. It is not covered by the template	no
Metro Manila : located in the administrative territorial entity : Philippines.	philippine' node is a neighor of both 'metro' and 'manila'	yes
Metro Manila : country : Philippines.	philippine' node is a neighor of both 'metro' and 'manila'	yes
Ninoy Aquino International Airport : located in the administrative territorial entity : Pasay City.	passay is a neighbor of nodes 'airport', 'ninoy' and 'aquino'	yes
Ninoy Aquino International Airport : country : Philippines.	philippine' is reachable from nodes 'airport', 'ninoy' and 'aquino'	yes
Asian Spirit : country : Philippines.	philippine' node is reachable for 'spirit' through 'airasia'	yes

Figure 5.18: Zest Airways Detection Summary

Summary

Although the validation has been performed on a single example (out of over 1200) in the DocRED dataset, the proposed method's capabilities can still be concluded.

The benchmark performance is evaluated for weakly supervised scenarios (Fig. 5.19). The weekly supervised scenario is a closer, but not the same, approach as the one proposed. This scenario provides additional training data containing the entities and expressing their exact relations. For example, "Ninoy Airport is located in Manila" would be extended with sentences such as "Ninoy Airport is the main airport in Manila". This approach resembles, but very roughly, semantic templates. Although it is impossible to compare the result directly, the Ign F1 score is at the same level as the F1 score for the proposed method.

Madal		Dev				Test		
Model	Ign F1	Ign AUC	F1	AUC	Ign F1	Ign AUC	F1	AUC
Supervised Setting								
CNN	37.99	31.47	43.45	39.41	36.44	30.44	42.33	38.98
LSTM	44.41	39.78	50.66	49.48	43.60	39.02	50.12	49.31
BiLSTM	45.12	40.93	50.95	50.27	44.73	40.40	51.06	50.43
Context-Aware	44.84	40.42	51.10	50.20	43.93	39.30	50.64	49.70
Weakly Supervised Setting								
CNN	26.35	14.18	42.75	38.01	25.40	13.46	42.02	36.86
LSTM	30.86	15.62	49.91	42.78	29.75	14.97	49.91	42.78
BiLSTM	32.05	16.50	51.72	44.42	29.96	15.50	49.82	42.90
Context-Aware	32.43	15.86	51.39	43.02	30.27	15.11	50.14	41.52

Figure 5.19: DocRED results. Ign F1 is calculated for dev/test datasets with removed duplicates [49]

The proposed solution is low-precision. It results from validating all relationship candidates available in the sample description, not just those defined in the training set. From a risk analysis perspective, reviewing all possible candidate entities is mandatory to identify a complete hazard flow. In contrast, the validation approach in DocRED assumes models are validated using examples in the dedicated validation set only. The DocRED precision scores are higher as they are calculated based on selected (fewer) cases vs. all possible pairs in the solution proposed.

The proposed solution implicitly considers the "no relation" case, another reason for low precision. "No relation" means that the entailment score is below the threshold and there is no link between the entities in the IRG. The DocRED dataset completely misses this case. As, in fact, "no relation" is a dominant relation, adding it would skew data distribution and severely impact the benchmark DocRED models' performance. Again, from a riskanalysis perspective, it is mandatory to explicitly review all candidate relations and assign "no relation" if there is not enough evidence supporting it.

5.3 EMARS Report Example

The Major Accident Reporting System ¹ (MARS and later renamed eMARS after going online) was first established by the EU's Seveso Directive 82/501/EEC in 1982 and has remained in place with subsequent revision to the Seveso Directive in effect today. eMARS contains reports of chemical accidents and near misses provided to the Major Accident Hazards Bureau (MAHB) of the European Commission's Joint Research Centre (JRC) from EU, EEA, OECD, and UNECE countries (under the TEIA Convention). Reporting an event into eMARS is compulsory for EU Member States when a Seveso establishment is involved, and the event meets the "major accident" criteria defined by Annex VI of the Seveso III Directive (201218/EU). The repository's goal is to facilitate the exchange of lessons learned from accidents and near misses involving dangerous substances to improve chemical accident prevention and mitigation of potential consequences.

The recorded reports are free-text descriptions of the events, unsuitable for comprehensive analysis without transformation to a network representation of hazard flow in the system domain. The proposed solution can extract the AVH triples from the description and transform the description of the event into the A-V-H graph.

Consider the example of ammonia leak (Fig. 5.20) which was provides with the following information:

Accident description

Release of ammonia, which intoxicated 11 persons (8 employees, of which 2 were intoxicated (injured) seriously and 3 fire-fighters). The release occurred due to a valve opened in error. The release occurred in an installation for the dilution of anhydrous ammonia into a 10% ammonia solution employed in order to limit (reduce) the

¹https://emars.jrc.ec.europa.eu/en/emars/content

corrosion during the distillation of petrol (crude oil). This release occurred on one (or two) 1/4 turn valves isolating the dissolving column from a tank of 7 cubic metres (m3) of capacity (containing 3,8 tonnes of liquefied ammonia under 8 bars of pressure at the time of the accident).

• Installation description

(anhydrous) ammonia dissolving (dilution) column employed to prepare the (10%) ammonia (solution) employed for its corrosion inhibiting characteristics during the (crude) oil refining process. The release occurred in correspondence of two flanges which were being made loose.

• Causes description

Probably the bad ergonomics of the place: during the unscrewing (loosening) of the bolts fixing the flanges, in a very cramped (small) space, one of the operators (workers) may have untimely (wrongly) opened one of the 1/4 turn valves isolating the unit (equipment) under maintenance from the units located upstream and containing the ammonia.

Consequences Description

3,8 tonnes of ammonia were released during 5 and a half hours (5,5 hours, 5,5h). The cloud stays (is confined) inside the establishment. 11 persons are affected or intoxicated (8 employees of which 2 were intoxicated (injured) seriously and 3 fire-fighters).

Emergency Response

Use of water and specifically for water curtains in order to abate the pollution (polluting substances - ammonia).

ent Profile		
Date/Time of Major	Occurrence	
Start Date and Time	24-09-1996 08:30	End Date and Time 24-09-1996 00:00
Accident ID	000758	
Accident Title	Release of ammonia from an e	roneous ly opened valve
Language	EN	Reported under EU Seveso I Directive
Event Type	Major Accident	Seveso II/III status - not known / not applicable -
Inductoin LA stivity		



Processing

All textual fields were altered for proper sentence formulation, i.e., repeated verbs were removed ("injured" as there is a verb "intoxicated"), and missing dots were added for correct sentence separation. All text was transformed into the SFG graph. The semantic pattern

"has a direct effect on" used was only. The cutoff threshold was 0.684703278413859. The relations above the threshold were loaded into the IRG graph and then aggregated into A-V-H.

Results

The event description contains a significant number of candidate relations (Tab. 5.2). However, 74% of them are irrelevant from a risk propagation perspective. Those accepted were loaded in the IRG where the PageRank measure [87] is calculated their "importance" in the overall semantic model of risk propagation (Tab. 5.3). The score is calculated regardless of the nodes' roles: (*Asset, Vulnerability* or *Hazard*). The high PageRank score of "ergonomic" seems correct, as the accident described resulted from a maintenance mistake due to a small space. The "ergonomic" node links with events, actors, and elements associated with the accident, such as "unscrewing," "operator," "turn," and "valve" (Fig. 5.21). The propagation of impact from these elements seems reasonable by reading the cause description.

The degree centrality of risk and vulnerability nodes (Tab 5.4) provides additional insides. In the A-V-H graph, "flange" (pol. kołnierz) is one of the most central *Risks* as it impacts indirectly 21 *Assets* through 15 *Vulnerabilities*. It is an essential *Vulnerability* as well (Fig. 5.22) as it connects and transmits the impact of other elements in the description. Such impact is confirmed by analysis of the "reasoning path" (Fig. 5.23). The path shows that it aggregates the impact of the following elements: the "valve" that is a part of the "equipment" isolating the ammonia. Then to the "operator" that was maintaining the "equipment". The "operator" was working in the "place" and making "turns" on the bolts in questionable "ergonomic" conditions, which eventually led to "flange" loosening that resulted in the "release".

The precision of the detection, however, is low, which means some relations accepted do not express a direct impact, and inspection of the "reasoning path" is required to confirm them. For example, "ammonia" is not a *Risk* for "ergonomic" (Fig. 5.24). The "reasoning path" does not confirm it, although the entailment score is above the threshold (Fig. 5.25).

Statistics	#
total number of candidate relations	1894
total number of accepted relations	483

Table 5.2: EMars report candidate and accepted relations counts

node	pagerank
valve	0.890
ergonomic	0.0850
release	0.0718
place	0.0691
equipment	0.0631
corrosion	0.0627
characteristic	0.0597
ammonia	0.0522
turn	0.0382
space	.0369
flange	0.0345

Table 5.3: The top IRG nodes given the pagerank



Figure 5.21: eMars example. The neighborhood of "ergonomic" node in the IRG graph



Figure 5.22: eMars example. The "flange" *Vulnerability* (Yellow node) in the A-V-H graph. The *Asset* nodes are blue and *Risk* nodes red

Path: [('valve', 'noun'), ('equipment', 'noun'), ('operator', 'noun'), ('place', 'noun'), ('turn', 'noun'), ('ergonomic', 'noun'), ('flange', 'noun'), ('release', 'noun')]

Edge: ('valve', 'noun') - ('equipment', 'noun') Properties: {'weight': 0.7624437212944031, 'frames': 'one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.', 'label': 'had a direct effect on'}

Edge: ('equipment', 'noun') - ('operator', 'noun') Properties: {'weight': 0.7229064702987671, 'frames': 'one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.probably the bad ergonomics of the place during the unscrewing of the bolts fixing the flanges in a very cramped small space one of the operators may untimely wrongly opened one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.', 'label': 'had a direct effect on'}

Edge: ('operator', 'noun') - ('place', 'noun') Properties: {'weight': 0.7751824855804443, 'frames': 'probably the bad ergonomics of the place during the unscrewing of the bolts fixing the flanges in a very cramped small space one of the operators may untimely wrongly opened one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.', 'label': 'had a direct effect on'}

Edge: <u>('place', 'noun') - ('turn', 'noun')</u> Properties: {'weight': 0.7402564287185669, 'frames': 'probably the bad ergonomics of the place during the unscrewing of the bolts fixing the flanges in a very cramped small space one of the operators may untimely wrongly opened one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.', 'label': 'had a direct effect on'}

Edge: ('turn', 'noun') - ('ergonomic', 'noun') Properties: {'weight': 0.8559473752975464, 'frames': 'one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.probably the bad ergonomics of the place during the unscrewing of the bolts fixing the flanges in a very cramped small space one of the operators may untimely wrongly opened one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.', 'label': 'had a direct effect on'}

Edge: <u>('ergonomic', 'noun') - ('flange', 'noun')</u> Properties: {'weight': 0.6959482431411743, 'frames': 'probably the bad ergonomics of the place during the unscrewing of the bolts fixing the flanges in a very cramped small space one of the operators may untimely wrongly opened one of the 1/4 turn valves isolating the unit equipment under maintenance from the units located upstream and containing the ammonia.the bolts fixing the flanges in a very cramped small space.', 'label': 'had a direct effect on'}

Edge: ('flange', 'noun') - ('release', 'noun') Properties: {'weight': 0.6918486952781677, 'frames': 'the release occurred in correspondence of two flanges which were being made loose.', 'label': 'had a direct effect on'}

Figure 5.23: eMars example. The path explaining "flange"'s impact on ammonia "release"



Figure 5.24: eMars example. "Ammonia"'s excessive *Vulnerability*. Blue nodes - *Assets*, Yellow - *Vulnerabilities*, red - *Risks*



Figure 5.25: eMars example. "Ammonia"->"ergonomic" reasoning path

Vulnerability	Risk
unit	capacity
ergonomic	tank
place	dilution
valve	flange
ammonia	correspondence
equipment	accident
unit	order
oil	column
petrol	ergonomic
solution	operator

Table 5.4: Top Vulnerability and Risk nodes sorted by their centrality measures in the A-V-H network

Summary

Given the accident description, the proposed solution correctly identified key elements, events, and actors in the infrastructure. The central element is the "valve" (Tab. 5.3), which is the main reason for the accident. The next two, "ergonomic" and "release", are correlated with it. The network risk model - the A-V-H graph combines all the elements of the description together. It represents the individual impact of all the elements on each other. It is possible to apply graph algorithms to focus the analysis on the elements with the highest impact and perform a more detailed analysis, such as one for "flange".

The "ammonia" example confirms the DocRED results as the precision of the detection is low. Low precision, however, can be attributed to the efficacy of the BART NLI classifier [82], its training set, and the efficacy of LLMs used to validate the results. As mentioned in the "Relevant NLP Processing Techniques," general NLI classifiers are trained using general NLI training sets, which may not cover the linguistic peculiarities of the specialized domain. On the other hand, performance of LLMs can be augmented by, for example, targeted examples that would improve prompting. Further work on precision improvement is required.

5.4 Impact of Large Langauge Models

As the introduction mentions, a few LLMs are explicitly trained in a dedicated domain, e.g., medical or legal. Unfortunately, the risk analysis domain does not have enough textual resources to perform such comprehensive training. On the other hand, LLMs e.q.: FLAN [60], [86], GPT3 [13] or Mistral [88] have achieved remarkable success lately establishing new reference results in a variety of general NLP tasks i.e. MMLU¹. LLMs have achieved such results following the latest in-context learning approaches such as instruction-tuning [60], chain-of-thoughts [89], or self-consistency [90]. The performance of the LLM is achieved because of good quality textual resources such as Wikipedia and thousands of books used in their training. However, there are differences in their training, which impacts their applicability to a specific language task, such as recognizing textual entailment.

The LLM is used to evaluate the entailment score threshold cutoff. The multiobjective optimization approach selects the cutoff for the entailment score so that the distance between the ratios of rejected and accepted relations is the largest. Therefore, it will select the point on the weight axis for which the distance between empirical distributions (ECDF) of rejected (blue curves) and accepted relations (orange curves) are as far as possible (Figs. 5.30, 5.32, 5.34). The statistics suggests that the vanilla models' performance (promoted in a general way, not specifically for the task) varies significantly. Therefore, a valid question arises: Which one should be chosen, and how does LLM performance relate to the RTE classification performed by the other classifier?

The experiment setup was as follows. The narrative described the Teton Dam collapse in 1975 (Fig. 5.29). The text was decomposed to the SFG representation. A single relation template "has a direct impact on" was used, and metapaths 1 and 2 were executed only to remove the effect of dialog consistency. The statistics (aggregated counts and empirical distribution) as a function of template entailment score - denoted as weight- were collected for LLMs (FLAN, LLAMA, Mistral) available on huggingface.com. Each model was prompted according to its interface with precisely the same prompt as others.

Together with the ECDFs, the distributions of LLM's *Yes / No* decisions are plotted against the entailment score for each model (Figs. 5.31, 5.33, 5.34, 5.36. The cutoff threshold for each model is provided in Fig. (5.37). It is impossible to assume that the threshold with a higher value is better. Therefore, the model selection has to rely on its performance for the task, not the threshold value itself.

¹https://paperswithcode.com/sota/multi-task-language-understanding-on-mmlu
The sample is large enough, and the collected statistics allow drawing the following assumption:

- the overall distribution of the number of relations per weight is exponential (Fig. 5.26)
- Most LLMS rejections shall be in low-scoring areas as most frames connecting nouns will not verbalize the relation template.
- the LLMS should agree in general with the RTE classifier therefore, the acceptance probability should be correlated with the weight and increase with it.(Fig. 5.27). In the perfect scenario, all rejected relations shall be close to 0 and accepted close to 1.

Therefore, the rudimentary criterion for selecting the performing model would be the difference *d* between the mean value of the score for acceptance and the mean value for rejection. Let w^0 denote the entitlement score (weight) associated with the LLM's rejection of the relations and w^1 the acceptance, then the distance

$$d = E[w^1] - E[w^0]$$

will denote the quality of LLMs decisions:

- if d>0, the LLM correlates with the RTE and properly discerns the rejections and acceptances. Solving the optimization task (Eq. 4.3) will establish the threshold,
- if d=0, the LLM is independent of RTE and cannot be used to estimate the threshold (Fig. 5.28),
- if d<0, the LLM is negatively correlated to the RTE and cannot be used to estimate the threshold.

The comparison of differences between the means for the models analyzed (Fig. 5.38) suggests using the mistral-v0.1 model. The acceptance threshold for this model has been calculated and is 0.273 (Fig. 5.37)



Figure 5.26: Histogram of relations given entitlement score (weight) for the exemplary narrative (Fig. 5.29)



Figure 5.27: Ideal histogram of exponential distribution of the number of relations given entailment score (weight) where class assignment correlates with entitlement score (weight)



Figure 5.28: Synthetic empirical distribution (left) and histogram (right) of the case where LLM classification is uncorrelated with the entailment score

Teton Canyon ends about six miles below the dam site, where the river flows onto the Snake River Plain. When the dam failed, the flood struck several communities immediately downstream, particularly Wilford at the terminus of the canyon, Sugar City, Salem, Hibbard, and Rexburg. Thousands of homes and businesses were destroyed. The small agricultural communities of Wilford and Sugar City were wiped from the river bank. Five of the eleven deaths attributed to the flood occurred in Wilford. The similar community of Teton, on the south bank of the river, is on a modest bench and was largely spared. One Teton resident was fishing on the river at the time of the dam failure and was drowned. An elderly woman living in the city of Teton died as a result of the evacuation. One estimate placed damage to Hibbard and Rexburg area, with a population of about 10,000, at 80% of existing structures. The Teton River flows through the industrial, commercial, and residential districts of north Rexburg. A significant reason for the massive damage in the community was the location of a lumber yard directly upstream. When the flood waters hit, thousands of logs were washed into town. Dozens of logs hit a bulk gasoline-storage tank a few hundred yards away. The gasoline ignited and sent flaming slicks adrift on the racing water. The force of the logs and cut lumber and the subsequent fires practically destroyed the city. The flood waters traveled west along the route of the Henrys Fork of the Snake River, around both sides of the Menan Buttes, significantly damaging the community of Roberts. The city of Idaho Falls, even further down on the flood plain, had time to prepare. At the older American Falls Dam downstream, engineers increased discharge by less than 5% before the flood arrived. That dam held and the flood was effectively over, but tens of thousands of acres of land near the river were stripped of fertile topsoil. The force of the failure destroyed the lower part of the Teton River, washing away riparian zones and reducing the canyon walls. This seriously damaged the stream's ecology and impacted the native Yellowstone cutthroat trout population. The force of the water and excessive sediment also damaged stream habitat in the Snake River and some tributaries, as far downstream as the Fort Hall bottoms. In August 1975, the region experienced an extreme flood, resulting in quantities of water falling that had not been considered during the construction of the dam. More than a year's worth of rain fell in only 24 hours, and the dam failed on August 8. Early on August 8, the dam was breached and 700 million cubic meters of floodwater was released, flooding communities and homes downstream. After this burst, a chain reaction began and the other 61 reservoirs located in the area collapsed—releasing another six billion cubic meters of floodwater. The water covered an area equal to 10 000 square kilometers. The failure to plan and account for extreme floodwaters resulted in the immediate death of 26 000 as a result of the water itself. 145 000 more people died as a result of epidemics and famine following the flood. The cause of the incident was reported to be unsafe vibrations coming from one of the turbines, which caused the turbine to break apart violently. Water that had been entering the turbine, flooded the turbine hall-flooding the room and levels below. The ceiling of the hall also broke apart from the impact from the turbine. At this point, power failed in the power station resulting in a blackout. Gasoline can cause fires. Steel gates to the water intake pipes of the turbines were manually closed and spillways were opened to prevent more damage.

Figure 5.29: Teton Dame Collapse Description



Figure 5.30: Emprical Distribution (ECDF) of LLama models per entailment score



Figure 5.31: Histograms of model's decisions on the relation per entailment score



Figure 5.32: Emprical Distribution (ECDF) of Mistral and chatGPT models per entailment score

Histogram of model decision on relations per entailment score Histogram of model decision on relations per entailment score



Figure 5.33: Histograms of model's decisions on the relation per entailment score



Histogram of model decision on relations per entailment score

Figure 5.34: Histograms of model's decisions on the relation per entailment score



Figure 5.35: Emprical Distribution (ECDF) of FLAN models per entailment score



Figure 5.36: Histograms of model's decisions on the relation per entailment score

model	relation	treshold
mistral-8x7B-Instruct-v0.1	has a direct impact on	0,273590982
chatGPT	has a direct impact on	0,213499486
mistral-7B-Instruct-v0.2	has a direct impact on	0,183974847
llama2_13b	has a direct impact on	0,158783749
llama2_7b	has a direct impact on	0,158696875
flan_t5_base	has a direct impact on	0,106075577

Figure 5.37: The cutoff thresholds for each LLMs

model	distance
mistral-8x7B-Instruct-v0,1	0,207326455
flan_t5_base	0,137173552
mistral-7B-Instruct-v0,2	0,106743931
llama2_13b	0,094051495
llama2_7b	0,092510897
chatGPT	0,059997983

Figure 5.38: Calculated distances between expected values for accepted and rejected relations per LLM

Summary

Surprisingly, models did not respond to the prompt consistently. For example, FLAN LLM showed a behavior in which models largest in terms of parameters (large and XXL) (Fig. 5.35) did not respond to the prompt at all, accepting the relations. This rudimentary analysis

shows differences between the models in how they encode the linguistic phenomena and, to some extent, explains the low precision of the overall detection quality. The discrepancies are not just between the models, which their training strategy can explain, but what is more surprising is that most of them do not show a correlation between the basic entailment score. For example, LLama models accept the majority of relations in the low entailment value regions (Fig. 5.30). For the future, vanilla prompting shall be more specific, and examples of entailment correlated with the domain shall be provided. Large Language Models need to be specifically verified for their applicability to the domain language.

5.5 Additional Examples

In other areas, non-linear systems exist, and it is essential to detect hazard propagation based on their description. In the financial scenario, for example, it is possible to model the impact of decisions given the economic phenomena they induce, i.e., inflation. Such an impact is challenging to predict due to the level of complexity of the market in terms of the interaction of dependent participants such as government, companies, and customers. Moreover, their connections are tight; therefore, the economic system fulfills the non-linear system definition. The effect propagates across all the participants, often with excessive feedback, which, if uncontrolled, can lead to critical phenomena such as hyperinflation or significant supply disruptions.

Another area that could benefit from the proposed methodology is drug development. The complexity of the human body is unquestionable. One identified challenge is preventing adverse drug effects, which aims to "predict the efficacy and toxicity of potential drug compounds" [91]. The proposed solution, relying on the description of the drug's mechanism of action under development, can help identify areas requiring special "attention" in clinical trials.

5.5.1 Financial Scenario

In this scenario, the description of price increase, i.e., inflation, is analyzed for potential impact on other market participants. The flow of inflation impact resembles the hazard propagation. The generated narrative (Fig. 5.39) describes the context surrounding the effect of inflation on the economy. Given the provided description, it is possible to explore the impact, such as the one presented for "inflation" for its closest neighbors (Fig. 5.40). The A-V-H graph aggregates the overall impact (Fig. 5.41), which can be used to further

Consumer price increase, often referred to as inflation, is the general increase in the prices of goods and services over time. It means that, on average, the cost of purchasing goods and services rises, leading to a decrease in the purchasing power of money. Demand-Pull Inflation occurs when demand for goods and services exceeds supply. When consumers have more money to spend or when there is increased demand due to factors like economic growth or government stimulus, businesses may raise prices to capitalize on the higher demand.Cost-Push Inflation happens when the cost of production for goods and services increases. Factors such as rising wages, increased raw material costs, or higher taxes can push up production costs, prompting businesses to pass these costs onto consumers through higher prices. Central banks control inflation by adjusting interest rates and the money supply. If a central bank increases the money supply excessively or keeps interest rates too low for too long, it can lead to higher inflation as more money chases the same amount of goods and services. Disruptions in the supply chain, such as natural disasters, geopolitical tensions, or pandemics, can lead to shortages of certain goods and services. When supply is limited and demand remains constant or increases, prices tend to rise. As prices rise, the same amount of money buys fewer goods and services, reducing the purchasing power of consumers' incomes.

Figure 5.39: Financial Example: Description of inflation

select elements of interest and inspect, for example, the "reasoning paths" (Figs. 5.42) 5.43).

5.5.2 Medical Scenario

In the medical scenario, it is important to identify an indirect impact of the drug on organs or functions that do not participate in the drug's mechanism of action. The sample text describes the role of the liver in the human body (Fig. 5.44). Assuming therapy impacts the liver directly, the goal is to identify potentially impacted elements or processes of the human body.

In this scenario, three relationship templates were used: "has effect on", "is a type of", "is a part of" and "is used by" (Fig. 5.45), the A-V-H graph aggregates the effect propagation into the combined network identifying the AVH triple (Fig. 5.46). The effect of the drug will be associated with "liver" in the *Vulnerability* role. The drug will impact the *Assets*, which are processes or organs associated with the "liver," i.e., "fibrinogen". To confirm the impact, the complete path should be evaluated. The path contains frames, the type of relation, and the entailment score (Fig. 5.47)



Figure 5.40: Financial Scenario: the "inflation" neighborhood in the IRG.



Figure 5.41: Financial Scenario: the A-V-H graph for "inflation" as Vulnerability of "business" Asset,Blue nodes - *Assets*, Yellow - *Vulnerabilities*, red - *Risks*

Path: ('government', 'noun') - ('inflation', 'noun') Edge: ('government', 'noun') - ('inflation', 'noun')
Properties:
{'weight': 0.4722381581771249,
'label': 'has a direct impact on',
'frames': 'when is increased demand due to factors like
<pre>economic growth or government stimulus businesses.demand pull inflation.'}</pre>

Figure 5.42: Financial Scenario: "government" impact on "inflation"

```
('money', 'noun') - ('business', 'noun'
Edge:
('money', 'noun') - ('business', 'noun')
Properties: {'weight': 0.8233859498310313,
'label': 'has a direct impact on',
'frames': 'the cost of purchasing goods and services leading
to a decrease in the purchasing power of money.businesses pass
these costs onto consumers through higher prices.'}
```

Figure 5.43: Financial Scenario: "money" impact on "business"

The liver is responsible for metabolizing nutrients from the food we eat, including carbohydrates, fats, and proteins. The liver converts glucose into glycogen for storage and releases it when blood sugar levels drop. Additionally, the liver metabolizes fats and produces bile, which aids in fat digestion and absorption in the intestines. Moreover, the liver detoxifies harmful substances, such as drugs, alcohol, and metabolic waste products, by breaking them down and facilitating their excretion from the body. The liver synthesizes albumin, which helps maintain blood volume and pressure. The liver synthesizes clotting factors, such as fibrinogen, prothrombin. The liver serves as a storage site for various nutrients and vitamins. It stores glycogen, which can be converted back into glucose when the body needs energy. The liver stores vitamins A, D, and B12, as well as iron and copper, which are essential for various metabolic processes. The liver stores excess glucose as glycogen during periods of high blood sugar and releases glucose into the bloodstream when blood sugar levels drop, ensuring that cells have a constant supply of energy. The liver produces bile, a greenish-yellow fluid that helps emulsify fats and facilitate their digestion and absorption in the small intestine. Bile is stored in the gallbladder and released into the small intestine when needed to aid in the digestion and absorption of dietary fats and fat-soluble vitamins. The liver plays a crucial role in the body's immune system by filtering and removing bacteria, viruses, and other pathogens from the bloodstream. The liver synthesizes various hormones and cholesterol necessary for maintaining hormonal balance and cell membrane integrity such as insulin-like growth factor 1 (IGF-1). Fibrinogen is a precursor to fibrin, which is the main protein component of blood clots.

Figure 5.44: Medical Example: description of the role of a liver



Figure 5.45: Liver example: the IRG graph



Figure 5.46: Liver example: the A-V-H graph. Blue nodes - Assets, Yellow - Vulnerabilities, red - Risks



Figure 5.47: Liver example: reasoning path for "liver" impact on "fibrinogen"

5.6 Summary

This chapter presented the solution's efficacy, focusing on analyzing and interpreting the use cases in alignment with the research objectives. It examined the exemplary narratives, highlighting relationship detection and anomalies observed during the analysis. The solution's functionality in risk analysis and applicability in modeling the flow of hazards in other domains have also been presented.

The quantitative analysis revealed several key issues with the approach that must be addressed. First, the low precision of risk-relationship detection is related to the efficacy of transformed-based classifiers and the current capabilities of large language models. Second, LLM prompting must be aligned with the domain to achieve "accept" and "reject" signals correlated with the entailment score for the threshold optimization task.

Chapter 6

Conclusions

The method described has demonstrated the following characteristics. First, it detects intra-sentence relations without training sets through a dialog consistency and verbalizing relationship pattern. Second, it shows that defining transitive relations applies to the risk analysis domain, as it is possible to construct the A-V-H, which is a risk-specific interaction graph. Third, we show that in the absence of training sets, a prompt-based classification using a language model can be used to provide a validation method. The proposed solution addresses the contextual entity classification problem and can be used to construct a comprehensive risk representation incrementally once new narratives are available.

Although there is a breadth of work on how current language models encode "knowledge" and how they can be used to extract it directly or validate the engineered hypothesis, limitations still hinder their direct applicability in various areas, including risk analysis. Hallucinations pose the most significant difficulty in their practical application. It seems that the relationship detection and classification task, quoted in this thesis as the relation verbalization, although being the most straightforward use case as it requires only properly defined prompt [48], [92] does not provide decisions as there are accepted relations that are significantly below the acceptance threshold. It seems reasonable to construct an ensemble of classifiers with prompted LLMs. Another limitation that pertains to LLMs' applicability is their context window fixed size. In in-context few-shot training, the context window would provide system descriptions. Therefore, descriptions that do not fit in must be split, resulting in possible lost relations. In addition, it is unclear how expanding the context window impacts the detection performance. It is also unclear how prompts would be constructed to solve the contextual representation classification (Risk-Asset-Vulnerability Dilemma). It seems that some sort of intermediate graph representation of text, similar to the SFG graph, which provides a method to limit the context (to a path of defined length) together with a specific graph traversing strategy, like metapaths, can help to fuse distant text fragments which would not fit into context if a description is provided directly. It seems that the proposed solution, to some extent, addresses the LLMs limitations mentioned and correctly uses LLMs as a *validation* instead of the main relationship *classification* or *detection* tool.

I want to continue research in the following directions. First, it is important to improve the precision of the detection. Second, I would like to focus on the capabilities of the SFG graph to check if it is possible to include more risk-oriented text operations. For example, to augment coreference resolution beyond simple preposition-noun or similar noun substitution, perform contextual substitution for the attributes of frames expressing the same event. For example, the sentences *"Electricity cutoff disables all electronic devices. In such situations, emergency UPSs provide backup power supply.", situation* refers to *electricity cutoff.* In the SFG graph, these sentences are unconnected as no nouns are in common. Therefore, either a specific connection or a complete replacement of *situation* with *electricity cutoff* should be performed to reduce the graph distance between *electronic devices* and *emergency UPS.*

Second, I would like to explore if it is possible to construct graphical representations of risk interaction other than A-V-H. I believe it is possible to model Event-Tree-Analysis by adding a dedicated set of templates. Combining them to detect how events co-occur and how the propagation of multiple hazards interacts in the system, we could create Fault-Tree-Analysis models [5]. Additionally, it might be interesting to verify if *path algebra* [93] can be applied to risk analysis.

Bibliography

- [1] A. Walczak, J. Napiórkowski, P. Adamczyk, and G. Kiryk, "Network model of risk analysis in the technical structures," in *21st International Conference on Circuits, Systems, Communications and Computers (CSCC 2017)*, 2017.
- [2] N. Leveson, "A new accident model for engineering safer systems," *Safety Science*, vol. 42, no. 4, pp. 237–270, 2004.
- [3] I. A. E. Agency, *Convention on nuclear safety*. International Atomic Energy Agency, 1994.
- [4] F. Crawley, A Guide to Hazard Identification Methods. Elsevier, 2020.
- [5] M. Rausand and S. Haugen, *Risk Assessment: Theory, Methods, and Applications*. WileySons, 2020.
- [6] C. Liu and S. Yang, "Using text mining to establish knowledge graph from accident/incident reports in risk assessment," *Expert Systems with Applications*, vol. 207, p. 117991, 2022.
- [7] F. Simone, S. M. Ansaldi, P. Agnello, and R. Patriarca, "Industrial safety management in the digital era: Constructing a knowledge graph from near misses," *Computers in Industry*, vol. 146, p. 103849, 2023.
- [8] V. Yadav and S. Bethard, "A survey on recent advances in named entity recognition from deep learning models," in *Proceedings of the 27th International Conference on Computational Linguistics* (E. M. Bender, L. Derczynski, and P. Isabelle, eds.), (Santa Fe, New Mexico, USA), pp. 2145–2158, Association for Computational Linguistics, Aug. 2018.
- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds.), vol. 30, Curran Associates, Inc., 2017.
- [10] C. D. Manning, "Human Language Understanding amp; Reasoning," *Daedalus*, vol. 151, pp. 127–138, 05 2022.

- [11] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. Chatterji, A. Chen, K. Creel, J. Q. Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, L. Fei-Fei, C. Finn, T. Gale, L. Gillespie, K. Goel, N. Goodman, S. Grossman, N. Guha, T. Hashimoto, P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. Krass, R. Krishna, R. Kuditipudi, A. Kumar, F. Ladhak, M. Lee, T. Lee, J. Leskovec, I. Levent, X. L. Li, X. Li, T. Ma, A. Malik, C. D. Manning, S. Mirchandani, E. Mitchell, Z. Munyikwa, S. Nair, A. Narayan, D. Narayanan, B. Newman, A. Nie, J. C. Niebles, H. Nilforoshan, J. Nyarko, G. Ogut, L. Orr, I. Papadimitriou, J. S. Park, C. Piech, E. Portelance, C. Potts, A. Raghunathan, R. Reich, H. Ren, F. Rong, Y. Roohani, C. Ruiz, J. Ryan, C. Ré, D. Sadigh, S. Sagawa, K. Santhanam, A. Shih, K. Srinivasan, A. Tamkin, R. Taori, A. W. Thomas, F. Tramèr, R. E. Wang, W. Wang, B. Wu, J. Wu, Y. Wu, S. M. Xie, M. Yasunaga, J. You, M. Zaharia, M. Zhang, T. Zhang, X. Zhang, Y. Zhang, L. Zheng, K. Zhou, and P. Liang, "On the opportunities and risks of foundation models," 2022.
- [12] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," 2020.
- [13] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems* (H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, eds.), vol. 33, pp. 1877–1901, Curran Associates, Inc., 2020.
- [14] J. Lai, W. Gan, J. Wu, Z. Qi, and P. S. Yu, "Large language models in law: A survey," 2023.
- [15] Y. Li, S. Wang, H. Ding, and H. Chen, "Large language models in finance: A survey," 2023.
- [16] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl, P. Payne, M. Seneviratne, P. Gamble, C. Kelly, N. Scharli, A. Chowdhery, P. Mansfield, B. A. y Arcas, D. Webster, G. S. Corrado, Y. Matias, K. Chou,

J. Gottweis, N. Tomasev, Y. Liu, A. Rajkomar, J. Barral, C. Semturs, A. Karthikesalingam, and V. Natarajan, "Large language models encode clinical knowledge," 2022.

- [17] A. G. M. Alam, "Semantic role labeling for knowledge graph extraction from text.," in *Progress in Artificial Intelligence*, 2021.
- [18] Y. Koreeda and C. Manning, "ContractNLI: A dataset for document-level natural language inference for contracts," in *Findings of the Association for Computational Linguistics: EMNLP 2021*, (Punta Cana, Dominican Republic), pp. 1907–1919, Association for Computational Linguistics, Nov. 2021.
- [19] C. Perrow, *Normal Accidents: Living with High Risk Technologies*. Princeton paperbacks, Princeton University Press, 1999.
- [20] J. J. Sammarco, "Operationalizing normal accident theory for safety-related computer systems," *Safety Science*, vol. 43, no. 9, pp. 697–714, 2005.
- [21] I. E. Comission, "Analysis techniques for system reliability procedure for failure mode and effects analysis, iec 60812," 2018.
- [22] T. A. Kletz, *Hazop Hazan: Identifying and Assessing Process Industry Hazards*. CRC Press, 1999.
- [23] T. R. Gruber, "A translation approach to portable ontology specifications," *Knowledge Acquisition*, vol. 5, no. 2, pp. 199–220, 1993.
- [24] R. Studer, V. Benjamins, and D. Fensel, "Knowledge engineering: Principles and methods," *Data Knowledge Engineering*, vol. 25, no. 1, pp. 161–197, 1998.
- [25] J. I. Single, J. Schmidt, and J. Denecke, "Knowledge acquisition from chemical accident databases using an ontology-based method and natural language processing," *Safety Science*, vol. 129, p. 104747, 2020.
- [26] P. Cimiano, "Ontology learning and population from text algorithms, evaluation and applications," 2006.
- [27] M. Hodkiewicz, J. W. Klüwer, C. Woods, T. Smoker, and E. Low, "An ontology for reasoning over engineering textual data stored in fmea spreadsheet tables," *Computers in Industry*, vol. 131, p. 103496, 2021.

- [28] I. O. for Standardization, Industrial Automation Systems and Integration-Integration of Life-cycle Data for Process Plants Including Oil and Gas Production Facilities. ISO, 2003.
- [29] P. Hughes, R. Robinson, M. Figueres-Esteban, and C. van Gulijk, "Extracting safety information from multi-lingual accident reports using an ontology-based approach," *Safety Science*, vol. 118, pp. 288–297, 2019.
- [30] S. M. Ansaldi, P. Agnello, A. Pirone, and M. R. Vallerotonda, "Near miss archive: A challenge to share knowledge among inspectors and improve seveso inspections," *Sustainability*, vol. 13, no. 15, 2021.
- [31] Z. Yin, L. Shi, Y. Yuan, X. Tan, and S. Xu, "A study on a knowledge graph construction method of safety reports for process industries," *Processes*, vol. 11, no. 1, 2023.
- [32] X. Zhao, H. Yan, and Y. Liu, "Hierarchical multi-label classification for fine-level event extraction from aviation accident reports," 2024.
- [33] R. Grishman and B. Sundheim, "Message Understanding Conference- 6: A brief history," in COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics, 1996.
- [34] L. Ramshaw and M. Marcus, "Text chunking using transformation-based learning," in *Third Workshop on Very Large Corpora*, 1995.
- [35] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference* of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) (J. Burstein, C. Doran, and T. Solorio, eds.), (Minneapolis, Minnesota), pp. 4171–4186, Association for Computational Linguistics, June 2019.
- [36] L. Ratinov and D. Roth, "Design challenges and misconceptions in named entity recognition," in *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL-2009)* (S. Stevenson and X. Carreras, eds.), (Boulder, Colorado), pp. 147–155, Association for Computational Linguistics, June 2009.
- [37] M. A. Hearst, "Automatic acquisition of hyponyms from large text corpora," in *COLING* 1992 Volume 2: The 14th International Conference on Computational Linguistics, 1992.

- [38] M. Honnibal and I. Montani, "spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing." To appear, 2017.
- [39] I. Hendrickx, S. N. Kim, Z. Kozareva, P. Nakov, D. Ó Séaghdha, S. Padó, M. Pennacchiotti,
 L. Romano, and S. Szpakowicz, "SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals," in *Proceedings of the 5th International Workshop on Semantic Evaluation* (K. Erk and C. Strapparava, eds.), (Uppsala, Sweden),
 pp. 33–38, Association for Computational Linguistics, July 2010.
- [40] R. Bunescu and R. Mooney, "A shortest path dependency kernel for relation extraction," in *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, (Vancouver, British Columbia, Canada), pp. 724–731, Association for Computational Linguistics, Oct. 2005.
- [41] N. Kambhatla, "Combining lexical, syntactic, and semantic features with maximum entropy models for information extraction," in *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, (Barcelona, Spain), pp. 178–181, Association for Computational Linguistics, July 2004.
- [42] B. Rosario and M. Hearst, "Classifying semantic relations in bioscience texts," in *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, (Barcelona, Spain), pp. 430–437, July 2004.
- [43] L. B. Soares, N. FitzGerald, J. Ling, and T. Kwiatkowski, "Matching the blanks: Distributional similarity for relation learning," 2019.
- [44] Y. Tian, G. Chen, Y. Song, and X. Wan, "Dependency-driven relation extraction with attentive graph convolutional networks," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, (Online), pp. 4458–4471, Association for Computational Linguistics, Aug. 2021.
- [45] T. T. Tran, P. Le, and S. Ananiadou, "Revisiting unsupervised relation extraction," in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, eds.), (Online), pp. 7498–7505, Association for Computational Linguistics, July 2020.

- [46] I. Dagan, O. Glickman, and B. Magnini, "The pascal recognising textual entailment challenge," in *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment* (J. Quiñonero-Candela, I. Dagan, B. Magnini, and F. d'Alché Buc, eds.), (Berlin, Heidelberg), pp. 177–190, Springer Berlin Heidelberg, 2006.
- [47] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, "A large annotated corpus for learning natural language inference," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, (Lisbon, Portugal), pp. 632–642, Association for Computational Linguistics, Sept. 2015.
- [48] O. Sainz, O. Lopez de Lacalle, G. Labaka, A. Barrena, and E. Agirre, "Label verbalization and entailment for effective zero and few-shot relation extraction," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, (Online and Punta Cana, Dominican Republic), pp. 1199–1212, Association for Computational Linguistics, Nov. 2021.
- [49] Y. Yao, D. Ye, P. Li, X. Han, Y. Lin, Z. Liu, Z. Liu, L. Huang, J. Zhou, and M. Sun, "DocRED: A large-scale document-level relation extraction dataset," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, (Florence, Italy), pp. 764–777, Association for Computational Linguistics, July 2019.
- [50] H. Tang, Y. Cao, Z. Zhang, J. Cao, F. Fang, S. Wang, and P. Yin, "Hin: Hierarchical inference network for document-level relation extraction," 2020.
- [51] J. Li, K. Xu, F. Li, H. Fei, Y. Ren, and D. Ji, "MRN: A locally and globally mention-based reasoning network for document-level relation extraction," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021* (C. Zong, F. Xia, W. Li, and R. Navigli, eds.), (Online), pp. 1359–1370, Association for Computational Linguistics, Aug. 2021.
- [52] Q. Huang, S. Zhu, Y. Feng, Y. Ye, Y. Lai, and D. Zhao, "Three sentences are all you need: Local path enhanced document relation extraction," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)* (C. Zong, F. Xia, W. Li, and R. Navigli, eds.), (Online), pp. 998–1004, Association for Computational Linguistics, Aug. 2021.

- [53] C. Quirk and H. Poon, "Distant supervision for relation extraction beyond the sentence boundary," in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers* (M. Lapata, P. Blunsom, and A. Koller, eds.), (Valencia, Spain), pp. 1171–1182, Association for Computational Linguistics, Apr. 2017.
- [54] F. Christopoulou, M. Miwa, and S. Ananiadou, "Connecting the dots: Document-level neural relation extraction with edge-oriented graphs," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (K. Inui, J. Jiang, V. Ng, and X. Wan, eds.), (Hong Kong, China), pp. 4925–4936, Association for Computational Linguistics, Nov. 2019.
- [55] D. Wang, W. Hu, E. Cao, and W. Sun, "Global-to-local neural networks for document-level relation extraction," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (B. Webber, T. Cohn, Y. He, and Y. Liu, eds.), (Online), pp. 3711–3721, Association for Computational Linguistics, Nov. 2020.
- [56] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling, "Modeling relational data with graph convolutional networks," 2017.
- [57] J. Li, Y. Wang, S. Zhang, and M. Zhang, "Rethinking document-level relation extraction: A reality check," 2023.
- [58] O. Levy, M. Seo, E. Choi, and L. Zettlemoyer, "Zero-shot relation extraction via reading comprehension," in *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)* (R. Levy and L. Specia, eds.), (Vancouver, Canada), pp. 333–342, Association for Computational Linguistics, Aug. 2017.
- [59] F. Petroni, T. Rocktäschel, S. Riedel, P. Lewis, A. Bakhtin, Y. Wu, and A. Miller, "Language models as knowledge bases?," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (K. Inui, J. Jiang, V. Ng, and X. Wan, eds.), (Hong Kong, China), pp. 2463–2473, Association for Computational Linguistics, Nov. 2019.
- [60] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," 2022.

- [61] Z. Wan, F. Cheng, Z. Mao, Q. Liu, H. Song, J. Li, and S. Kurohashi, "GPT-RE: In-context learning for relation extraction using large language models," in *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing (H. Bouamor, J. Pino, and K. Bali, eds.), (Singapore), pp. 3534–3547, Association for Computational Linguistics, Dec. 2023.
- [62] G. Li, P. Wang, and W. Ke, "Revisiting large language models as zero-shot relation extractors," 2023.
- [63] I. Dagan, D. Roth, M. Sammons, and F. M. Zanzotto, *Recognizing Textual Entailment: Models and Applications*. Synthesis Lectures on Human Language Technologies, Morgan & Claypool Publishers, 2013.
- [64] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," 2019.
- [65] A. Williams, N. Nangia, and S. R. Bowman, "A broad-coverage challenge corpus for sentence understanding through inference," 2018.
- [66] C. Fillmore, "Scenes-and-frames semantics," Linguistic Structures Processing, 1977.
- [67] M. Palmer, D. Gildea, and P. Kingsbury, "The Proposition Bank: An annotated corpus of semantic roles," *Computational Linguistics*, vol. 31, no. 1, pp. 71–106, 2005.
- [68] V. P. A. Gangemi, "Towards a pattern science for the semantic web.," in *Semantic Web Journal*, IOS Press, 2010.
- [69] D. Davidson, "The logical form of action sentences," in *The Logic of Decision and Action* (N. Rescher, ed.), pp. 81–95, University of Pittsburgh Press, 1967.
- [70] D. D.Gildea, "Automatic labeling of semantic roles," in *Association for Computational Linguistics*, ACL, 2000.
- [71] J. D.Jurafsky, Speech and Language Processing. Stanford: DRAFT, 2nd ed., 2022.
- [72] P. Shi and J. Lin, "Simple bert models for relation extraction and semantic role labeling," 2019.
- [73] D. Hendrycks, X. Liu, E. Wallace, A. Dziedzic, R. Krishnan, and D. Song, "Pretrained transformers improve out-of-distribution robustness," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (D. Jurafsky, J. Chai,

N. Schluter, and J. Tetreault, eds.), (Online), pp. 2744–2751, Association for Computational Linguistics, July 2020.

- [74] D. Yogatama, C. de Masson d'Autume, J. Connor, T. Kocisky, M. Chrzanowski, L. Kong, A. Lazaridou, W. Ling, L. Yu, C. Dyer, and P. Blunsom, "Learning and evaluating general linguistic intelligence," 2019.
- [75] T. McCoy, E. Pavlick, and T. Linzen, "Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (A. Korhonen, D. Traum, and L. Màrquez, eds.), (Florence, Italy), pp. 3428–3448, Association for Computational Linguistics, July 2019.
- [76] S. Gururangan, S. Swayamdipta, O. Levy, R. Schwartz, S. Bowman, and N. A. Smith, "Annotation artifacts in natural language inference data," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)* (M. Walker, H. Ji, and A. Stent, eds.), (New Orleans, Louisiana), pp. 107–112, Association for Computational Linguistics, June 2018.
- [77] Adamczyk, Piotr, Kiryk, Grzegorz, Napiórkowski, Jarosław, and Walczak, Andrzej, "Network model of security system," *MATEC Web Conf.*, vol. 76, p. 02002, 2016.
- [78] C. J. Fillmore, "The case for case," in *Universals in Linguistic Theory* (E. Bach and R. T. Harms, eds.), pp. 0–88, New York: Holt, Rinehart and Winston, 1968.
- [79] Y. Dong, N. V. Chawla, and A. Swami, "Metapath2vec: Scalable representation learning for heterogeneous networks," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, (New York, NY, USA), p. 135–144, Association for Computing Machinery, 2017.
- [80] B. J. Grosz, A. K. Joshi, and S. Weinstein, "Centering: A framework for modeling the local coherence of discourse," *Computational Linguistics*, vol. 21, no. 2, pp. 203–225, 1995.
- [81] R. Barzilay and M. Lapata, "Modeling local coherence: An entity-based approach," in Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05) (K. Knight, H. T. Ng, and K. Oflazer, eds.), (Ann Arbor, Michigan), pp. 141–148, Association for Computational Linguistics, June 2005.

- [82] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, (Online), pp. 7871–7880, Association for Computational Linguistics, July 2020.
- [83] "Chatgpt prompt engineering." http://https://www.promptingguide.ai/ models/chatgpt. Accessed: 2023-09-10.
- [84] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Comput. Surv.*, vol. 55, jan 2023.
- [85] A. P. Bradley, "The use of the area under the roc curve in the evaluation of machine learning algorithms," *Pattern Recognition*, vol. 30, no. 7, pp. 1145–1159, 1997.
- [86] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023.
- [87] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank Citation Ranking: Bringing Order to the Web," tech. rep., Stanford Digital Library Technologies Project, 1998.
- [88] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas,
 F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock,
 T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed, "Mistral 7b," 2023.
- [89] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2023.
- [90] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," 2023.
- [91] A. Blanco-González, A. Cabezón, A. Seco-González, D. Conde-Torres, P. Antelo-Riveiro, Piñeiro, and R. Garcia-Fandino, "The role of ai in drug discovery: Challenges, opportunities, and strategies," *Pharmaceuticals*, vol. 16, no. 6, 2023.
- [92] T. Schick and H. Schütze, "Exploiting cloze-questions for few-shot text classification and natural language inference," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume* (P. Merlo,

J. Tiedemann, and R. Tsarfaty, eds.), (Online), pp. 255–269, Association for Computational Linguistics, Apr. 2021.

[93] R. Bakker, *Knowledge Graphs : Representation and Structuring of Scientific Knowledge*.PhD thesis, University of Twente, Enschede, The Netherlands, 1987.

Appendix

.1 Computation Complexities Calculatation

.1.1 Naive Relation Validation

Each context c has k nouns on average. In the worst-case scenario, to validate the semantic relation between any pair of nouns in the document, we have to visit all contexts and all nouns in them. The validation is performed as follows:

- 1. we have to find a noun in one of the c contexts. In the worst-case scenario, we must scan c contexts and k nouns in each to find the noun of interest.
- 2. within the first context, we have to validate the relationship between the noun of interest and the remaining k-1 nouns,
- 3. then we need to find a connection between the current and next contexts. Therefore, we must scan the remaining c-1 context for shared nouns to make a connection. In the worst-case scenario, there are c-1 contexts left and for each, we need to scan k nouns
- 4. within the following context, we need to evaluate the relationship between the connecting noun and the remaining k-1 nouns
- 5. we repeat the context evaluation until the last context.

The formula estimating the number of comparisons depends on the number of contexts



Figure 1: Naive Validation: Each context has k nouns on average; there are context in total in the narrative

and average number of nouns:

$$f(k,c) = c \ast k \ast (k-1) \ast (c-1) \ast k \ast (k-1) \ast (c-2) \ast k \ast k(1) \ast \ldots (c-c-1) \ast k \ast (k-1)$$

$$f(k,c) = c!k^c(k-1)^c$$