# A conceptual framework for the automated analysis of railway accident reports

Wei-Ting Hong[1], Geoffrey Clifton[1], John D Nelson[1]

[1] Institute of Transport and Logistics Studies, The University of Sydney Business School

Email for correspondence: wei-ting.hong@sydney.edu.au

## 1. Introduction

The analysis of railway accident reports is fundamental to understanding the mechanism of hazards and observing how railway industries in different countries manage risks in the railway system. Although considerable prior work has attempted to increase the number of railway accident cases investigated, many of them suffer from the need to process a large amount of textual data. To overcome this obstacle, a conceptual framework for the automated analysis of railway accident reports is proposed. Only open-sourced data and coding packages are applied to building the model and the input data is not constrained by the length of reports and the variation in the use of English language. Two main Natural Language Processing (NLP) topic models are utilised, BERTopic and the Structural Topic Model (STM). A demonstration to illustrate the application of the framework proposed to the real-world data using railway accident reports published by the Rail Accident Investigation Branch (RAIB) in the United Kingdom, Australian Transport Safety Bureau (ATSB), National Transportation Safety Board (NTSB), USA and Transportation Safety Board of Canada (TSB) is presented. The result shows strong potential to automatically extract hazards and help stakeholders learn across jurisdictions. Future research could incorporate data from additional jurisdictions and the framework could be applied to road, aviation or maritime accidents.

## 2. Literature context

A growing number of studies place emphasis on railway safety through application of state-of-the-art NLP techniques, resulting in a significant opportunity to eliminate the restrictions on the analysis of big textual safety-related data. Some studies have indicated the possibility to classify railway accidents based on the features of original accident records (Hadj-Mabrouk, 2020). However, the determination of knowledge about critical hazards triggering the accident is still reliant on human determination, which is time-consuming and labor-intensive (Kim &Yoon, 2013; Zhou &Lei, 2018). Such limitations prevent updating of the model after new railway accident cases become available and hinder the development of railway infrastructure technology. Furthermore, railway safety-related frameworks published in the literature are seldom shared, making the data un-reusable and restricting further research.

## 3. Proposed conceptual framework

To automate the analysis of railway accident reports, topic modelling methods are leveraged to extract the relationship between topics and documents by different features such as the probability of occurrence of words and high dimensional word embeddings. Topic models assume that a document contains a collection of underlying themes and the distribution of words in the document over the whole corpus might derive topics representing these underlying themes. A set of keywords is identified to reflect underlying topics and their trend and statistics are derived for further methodological and practical applications (Blei and Mcauliffe, 2007).

A topic model can be trained in several ways, including supervised learning, semi-supervised learning and unsupervised learning. Unsupervised learning approaches are selected to build the topic model to ensure highly automated analysis and avoid human intervention. Several package-oriented programming models have been developed based on these packages and result in significant improvements in performance in the topic modelling contexts. The Structural Topic Model (STM) (Kwayu et al., 2021; Li et al., 2011) and BERTopic model (Grootendorst, 2022) have been commonly applied to NLP tasks due to the high performance achieved and the flexibility of estimate effect analysis (Labusch &Neudecker, 2020; Ly et al., 2020).

## 3.1. BERTopic

The BERTopic [1] is an open access topic model adopting the Bidirectional Encoder Representations from Transformers (BERT) pre-trained language model (Devlin et al., 2018) to retrieve high-dimension vectors of texts for clustering. For implementation, topics are generated through three steps: text vectorisation with a pre-trained language model, dimension reduction for optimising the modelling process, and topic representations with custom class-based TF-IDF (c-TF-IDF). The c-TF-IDF is an advanced method for converting original text into a series of representative numbers (which is also known as word embedding). In contrast to traditional approaches, the c-TF-IDF takes the semantic relationships between words into account, increasing the interpretability and accuracy of the outcomes. More mathematical details and theorems can be found in Devlin et al. (2018).

## 3.2. Structural Topic Model

The Structural Topic Model (STM)[2] is another open access and unsupervised learning-based probabilistic topic modelling method derived from the Latent Dirichlet Allocation (LDA). The LDA is a generative statistical model that classifies documents based on the observations of each individual word collected in the documents and assumes that the topic of each document is derived from the aggregation of the words in that document. The STM is developed on the same statistical basis as the LDA in addition to allowing correlations of external factors among topics. The main difference lies in the pre-generalised linear models derived from the nature of the data used while estimating parameters. More mathematical details and theorems can be found in Roberts et al. (2013).

The STM is more suitable than the LDA for analysis of railway accident reports because critical covariates are usually disclosed and discussed in reports, such as the occurrence of time and the relevant modes of rail transport (such as light, suburban, and heavy rail) and organisations. These critical covariates can offer valuable insights for better understanding the nature and prevalence of railway accidents across time. For instance, the STM may provide the difference in how the platform-train interface incidents occur on the light rail system and other modes of the rail transport system. The trend of the causes of an accident may also be revealed by supplementing the occurrence of time as an additional covariate in STM temporal analysis.

## 3.3. The data requirements

To demonstrate the application of the proposed models, railway accident reports published by independent railway accident investigation bodies from the UK, the USA, Canada and Australia are used. Railway accident reports compiled by independent railway accident investigation organisations are regulated by a national legislative framework and provide unbiased and blame-free details for promoting a railway safety culture. Additionally, only investigating

---

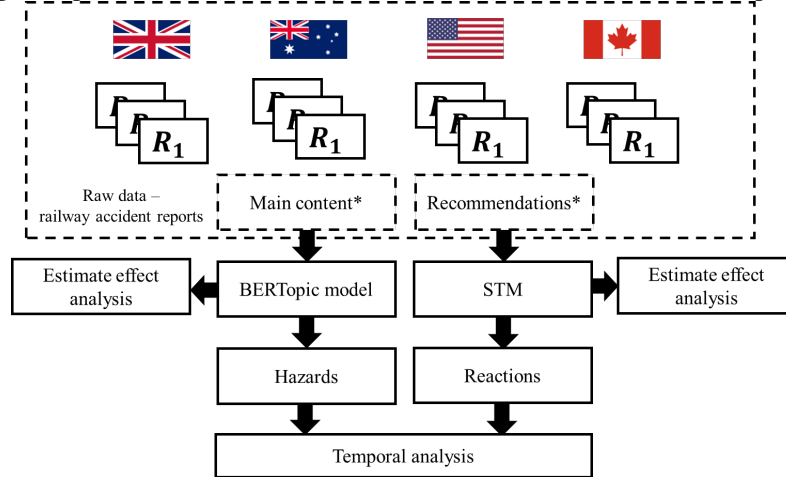[1] More details about BERTopic can be found in the following GitHub: https://github.com/MaartenGr/BERTopic
[2] More details about STM can be found in the following page: https://www.structuraltopicmodel.com/

bodies which have published over 100 reports with concrete recommendations and written in formal English are considered to ensure the performance of the model. Despite the differences in writing styles and terminology used, all reports consist of the summary of the accident, the analysis, the investigation, key findings, conclusions and recommendations (if applicable). In this study scanned files are removed due to the technical difficulties of recognising the text.

Each section of a railway accident report illustrates the accident from different perspectives. In addition, a mixture of critical information is also outlined, including causal factors, underlying factors, contributing factors, and identified hazards. All of these indicate the fact that a railway accident report contains a wide range of heterogeneous information that might not be fully captured by the document-level analysis. Therefore, BERTopic is appropriate for analysis at the sentence-level to handle the heterogeneity. On the other hand, recommendations have a strong semantic homogeneity of descriptions. In this case, the STM might be more applicable because the occurrence of words is more meaningful than the semantic context information.

The overview of the data flow and analysis procedures is illustrated in Figure 1. The main content of railway accident reports is used to extract potential hazards, whereas the recommendations made in railway accident reports are extracted and analysed separately for the purposes of understanding how each investigation body reacts to risks identified. Both outcomes are further extended for estimate effect analysis and temporal analysis, providing additional insights into solutions during different periods of time. To sum up, this work contributes to the railway safety context by offering the opportunity to look at a large volume of railway accidents from multiple perspectives and allowing end-users to have a comprehensive view of hazards across jurisdictions and time.

Figure 1: The overview of the data flow and analysis procedures for automated analysis of railway accident reports. Note that railway accident reports are divided into main content (including all descriptions such as summary, investigation process and conclusion) and recommendations for different purposes of analysis.
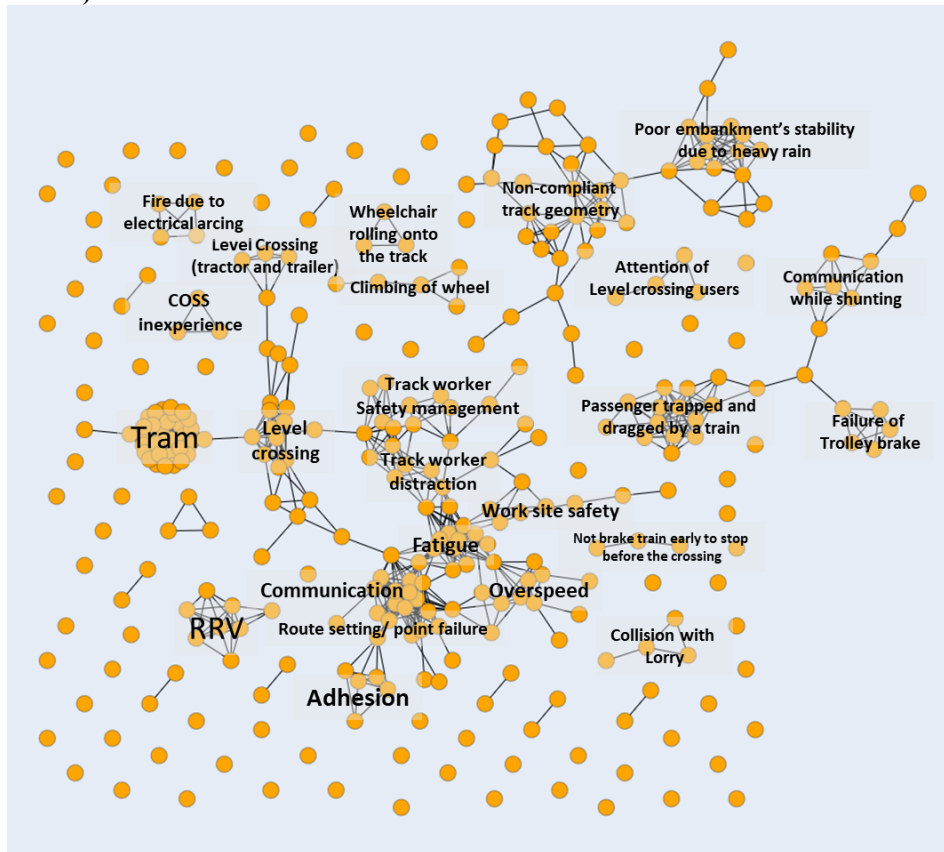


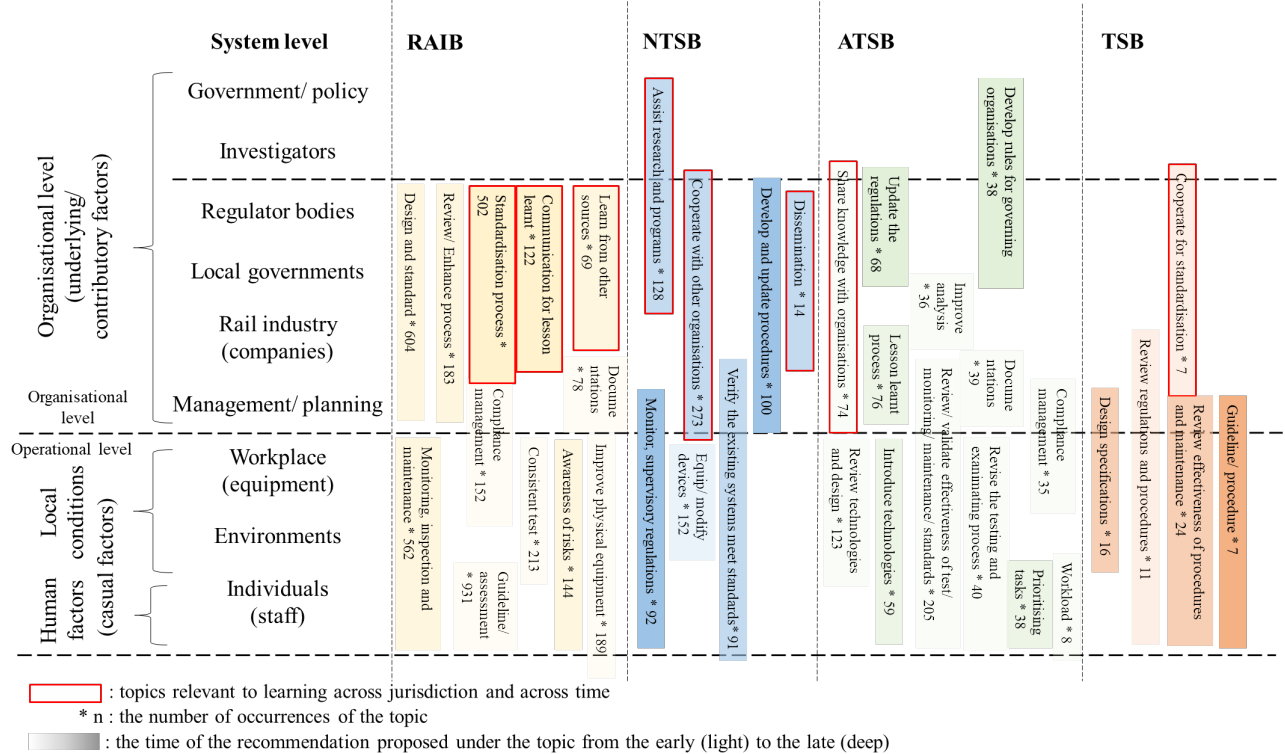## 4. A demonstration of practical implementation

The outcome of topic models only shows topics extracted and representative keywords of each topic. Nevertheless, the relationship between each topic cannot be revealed, hindering users from understanding the mechanisms of railway accidents. Therefore, the result of the BERTopic model is further extended by adding additional processes to address this issue. Firstly, the distribution of the number of sentences over each topic on documents is extracted and condensed to a topic-document matrix. Secondly, we assume that the distribution of each topic over documents is the projection of the extent to which this topic influences each railway accident. Multiple similar distributions indicate that these topics constitute a specific group of railway accidents with similar features. Therefore, the cosine similarity approach is applied to

identify the similarity of distributions (Cheng et al., 2009; Qurashi et al., 2020). A topic-topic similarity matrix can be generated with each element between 0 and 1. The larger similarity score indicates that sentences under both topics are commonly used in the same group of documents. Next, a distribution of topics, including the relationship, can be mapped by setting a threshold for the similarity score, linking each topic and forming a series of clusters representing various hazards. The threshold for the similarity might be determined based on the nature of input data and analysis purposes. Figure 2 shows the distribution of the relationship between hazards identified in the RAIB dataset (covered period: 2005-2019). Each orange dot represents a topic identified by the BERTopic and the link refers to the similarity score of two topics that is larger than the threshold. The name assigned to each cluster is based on the inference of keywords of linked topics. According to this result, several major concepts of hazards and their underlying causal relations have been revealed. Each group is dominated by a specific hazard along with several supplementary conditions. For example, there is a main topic "Communication" close to the central area and another communication-related topic "Communication while shunting" on the right-hand side. Despite the same keywords, the concept of "Communication" in the first case is the major topic, whereas in the second case it is the supplementary topic. Therefore, there is no connection shown in the figure due to the weak relationship.

**Figure 2: The distribution of the relationship between hazards identified in the RAIB dataset (covered period: 2005-2019)**



To depict the role each recommendation plays in the railway system and proposes, topics extracted from STM are fitted at multiple system levels and describe the trend of each type of recommendation over different countries. Figure 3 shows the distribution of recommendation topics at each socio-technical level. The light colour refers to the early recommendations and vice versa. The railway system is divided into the organisational level and operational level, representing how the socio-technical system works in the railway industry.

**Figure 3: The recommendations made by investigators at each socio-technical level[3]**



: topics relevant to learning across jurisdiction and across time
* n : the number of occurrences of the topic
: the time of the recommendation proposed under the topic from the early (light) to the late (deep)

Overall, Figure 3 maps out how investigators in different countries address identified hazards and lead the railway industry of each jurisdiction to improve railway safety. Common recommendations at the operational level are procedures of maintenance and inspection, consistency of testing processes, introducing state-of-the-art equipment, and reviewing existing designs and technologies. On the other hand, recommendations at the organisational level popularly proposed are process standardisation, co-operation with other organisations and dissemination of railway safety knowledge.

A considerable number of recommendations related to learning across jurisdictions and times (outlined in red) are proposed by RAIB and NTSB, implying a solid promotion of sharing knowledge with other research organisations. The trend continues nowadays along with recommendations relating to the dissemination of railway safety knowledge. It is also worth noting that NTSB consistently tends to propose precise but interfering recommendations, such as verifying existing systems and assisting research and programs. On the other hand, several recommendations made by ATSB and TSB indicate detailed instructions at the operational level, such as the prioritisation of tasks, the management of workload and validation of the effectiveness of existing standards but lack recommendations relevant to the learning behaviour. Despite no direct evidence of affecting railway safety, insufficient learning across jurisdictions and time might lead to a poor railway safety culture due to passive reactions to risks found in other countries.

# 5. Conclusion

This study proposes a conceptual framework for the automated analysis of railway safety reports. Topic models BERTopic and STM are utilised, and original outcomes are further extended for better interpretation. The result outlines the distribution of hazards and indicates how investigators react to risks identified. The learning behaviour of investigators in each

---

[3] An original size of Figure 3 is available on the following website: The recommendations made by investigators at each socio-technical level.png

jurisdiction is also revealed, suggesting insufficient learning across jurisdictions. The conceptual framework proposed provides a potential solution to reduce human effort required by automatically extracting critical hazards from a collection of railway accident reports published in different jurisdictions. The framework also enables the possibility to share knowledge by synthesising written knowledge and recommendations made across jurisdictions. Despite rich findings, the outcome of models still requires adequate human interpretation, especially in interpreting keywords extracted and assigning topic' names. The framework proposed can be applied to other textual data for further insights, such as road accidents, aviation safety and marine operations.

# 6. References

Blei, D. M., &Mcauliffe, J. D. (2007). Supervised Topic Models. *Advances in Neural Information Processing Systems*, *20*. www.digg.com

Cheng, M. Y., Tsai, H. C., &Chiu, Y. H. (2009). Fuzzy case-based reasoning for coping with construction disputes. *Expert Systems with Applications*, *36*(2 PART 2), 4106–4113. https://doi.org/10.1016/j.eswa.2008.03.025

Devlin, J., Chang, M.-W., Lee, K., &Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, *1*, 4171–4186. http://arxiv.org/abs/1810.04805

Grootendorst, M. (2022). *BERTopic: Neural topic modeling with a class-based TF-IDF procedure*. https://doi.org/10.48550/arxiv.2203.05794

Hadj-Mabrouk, H. (2020). Analysis and prediction of railway accident risks using machine learning. *AIMS Electronics and Electrical Engineering*, *4*(1), 19–46. https://doi.org/10.3934/electreng.2020.1.19

Kim, D. S., &Yoon, W. C. (2013). An accident causation model for the railway industry: Application of the model to 80 rail accident investigation reports from the UK. *Safety Science*, *60*, 57–68. https://doi.org/10.1016/j.ssci.2013.06.010

Kwayu, K. M., Kwigizile, V., Lee, K., &Oh, J. S. (2021). Discovering latent themes in traffic fatal crash narratives using text mining analytics and network topology. *Accident Analysis & Prevention*, *150*, 105899. https://doi.org/10.1016/J.AAP.2020.105899

Labusch, K., &Neudecker, C. (2020). Named Entity Disambiguation and Linking on Historic Newspaper OCR with BERT. *Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum*. https://qurator.ai

Li, S.Bin, Wu, J. J., Gao, Z. Y., Lin, Y., &Fu, B. B. (2011). The analysis of traffic congestion and dynamic propagation properties based on complex network. *Wuli Xuebao/Acta Physica Sinica*, *60*(5), 050701–050701. https://doi.org/10.7498/aps.60.050701

Ly, A., Uthayasooriyar, B., &Wang, T. (2020). A survey on natural language processing (nlp) and applications in insurance. *ArXiv*. http://arxiv.org/abs/2010.00462

Qurashi, A. W., Holmes, V., &Johnson, A. P. (2020). Document Processing: Methods for Semantic Text Similarity Analysis. *INISTA 2020 - 2020 International Conference on INnovations in Intelligent SysTems and Applications, Proceedings*. https://doi.org/10.1109/INISTA49547.2020.9194665

Roberts, M. E., Stewart, B. M., Tingley, D., &Airoldi, E. M. (2013). The structural topic model and applied social science. *Adv. Neural Inf. Process. Syst. Workshop Top. Models: Comput. Appl. Eval.*, 1–20.

Zhou, J. L., &Lei, Y. (2018). Paths between latent and active errors: Analysis of 407 railway accidents/incidents' causes in China. *Safety Science*, *110*(November 2017), 47–58. https://doi.org/10.1016/j.ssci.2017.12.027