# Executive Summary
Large Language Models and Topic Modelling for Textual
Information Extraction and Analysis

**Marcos Tidball**
GitHub: @zysymu
LinkedIn: in/zysymu

## Key Findings

- "Big narratives" (a custom prompt in human-readable format that uses a combination of tabular data and narratives) provide LLMs a more complete understanding of what happened.

- Text2Text Generation transformers that combine "big narratives" with specific questions allow for the automatic extraction of information present only on the narratives, such as precipitating events, activity involved and more complete diagnoses.

- Topic modelling can be used to create categories from information extracted from narratives.

## Summary of Approach

### Data Sources

The only data used is the Primary Data. I "translate" the most frequent technical terms and expressions present in the narratives into common terms (Fig. 1). This step was fundamental for optimal results when using pretrained Text2Text Generation LLMs.

### Methodology

Inspired by TabLLM, I create custom prompts ("big narratives") that combine patient information contained on tabular variables with the translated narrative (Fig. 2). I combine each big narrative with 6 questions (Fig. 3) and use the `google/flan-t5-base` model for Text2Text Generation; this allows the LLM to extract specific information about the falls based on the questions.

To analyze the answers, I trained one Latent Dirichlet Allocation (LDA) model per question to model the topics in the answers. This allows answers to be categorized into topics with similar subjects, which, together with the 6 different question categories provide a very granular way to analyze the data. I preprocessed the answers with: word-level tokenization, stop words removal, creation of bigrams and trigrams and lemmatization. However, even after hyperparameter optimization, some topics are still not very clear. This is probably due to the loss of positional and textual information caused by the preprocessing pipeline. I believe that using a more robust topic modelling algorithm such as BERTopic would provide better results.

### Evaluation

For the information extraction part, I tried different Question Answering models, but due to their extractive nature (answers could only be text present in the prompt), sparse narratives were problematic. I also tried using just the narratives available on the base dataset with and without the translation step, but the translated big narratives allowed for more complete and accurate answers, at least in a sample of 50 randomly selected narratives that I manually analyzed.

The LDA models were tuned only for coherence in mind, with their performance being evaluated via the coherence score, visualizations and by checking the most representative answers for each

topic (Fig. 4). Initial analyses and plots that demonstrate how these topics can be used are shown in Fig. 5 and Fig. 6.

## Visualizations



Figure 1: Translation dictionary for technical terms.



Figure 2: Example of automated big narrative generated for `cpsc_case_number` = 230217170. Original narrative is: "81 YOM FELL OUT OF BED. DX: LEFT FEMUR FRACTURE."



Figure 3: Questions that, combined with the big narratives, form the prompts fed to the LLM.



Figure 4: Most representative translated narratives (highest score) for topic 0 of the `cause` question. The keywords for this topic are: "slip, rugs_carpet, ceilings_wall, bathtubs_shower, frame, door_sill, shoe, walker, foot, catch".
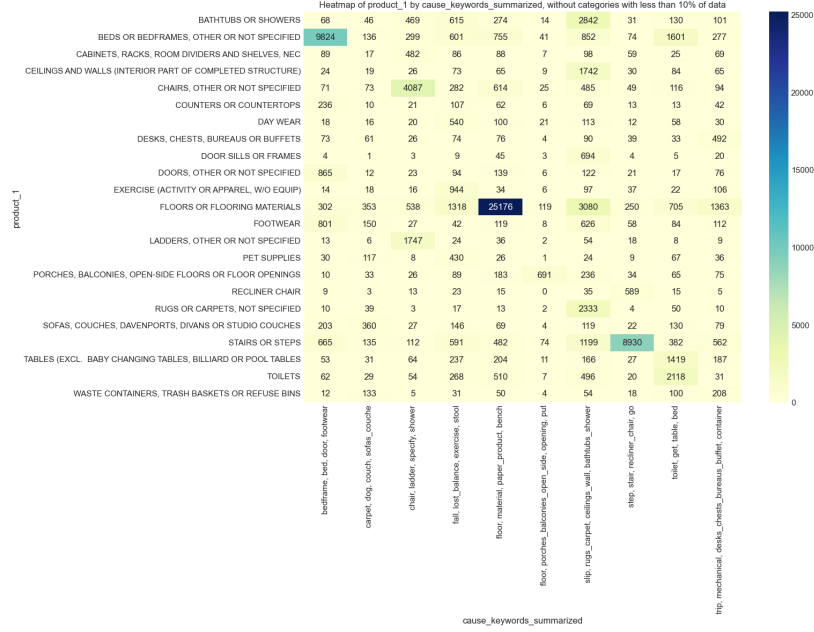
Figure 5: Heatmap of `product_1` by `cause` topics (with each topic's keywords being displayed for better explainability) displaying only categories with more than 10% of total data for better readability. We can notice that there is a good match between the `product_1` categorical variable and the more specific `cause` topics, which give us more information about what happened.
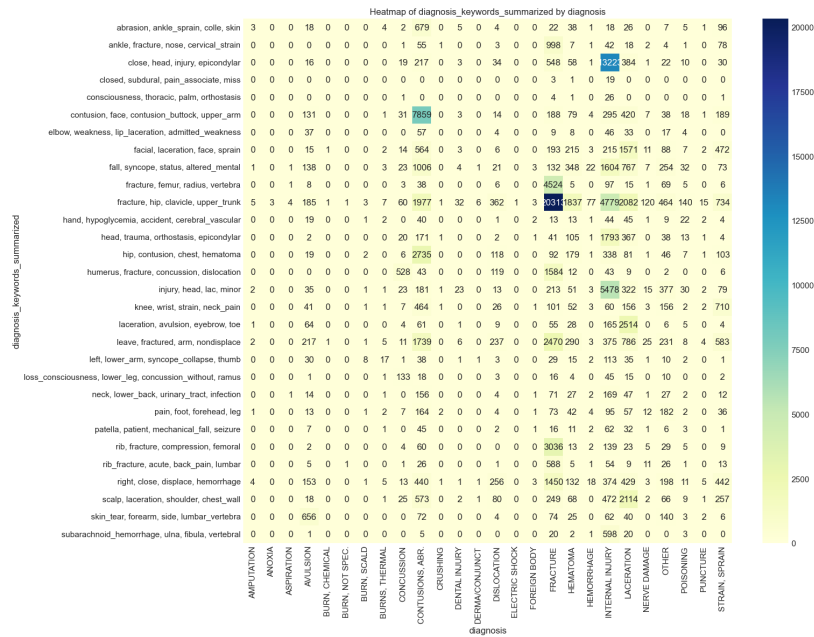


Figure 6: Heatmap of `diagnosis` by `diagnosis` topics (with each topic's keywords being displayed for better explainability). Visualizations like this can be used to evaluate if the topics make sense. This heatmap demonstrates that the information extraction pipeline that was developed can aid in obtaining more specific details about what happened in the fall events.