# Building Safer Sites: A Large-Scale Multi-Level Dataset for Construction Safety Research

Zhenhui Ou*
Arizona State University
Tempe, AZ, USA
zhenhuio@asu.edu

Dawei Li*
Arizona State University
Tempe, AZ, USA
daweili5@asu.edu

Zhen Tan
Arizona State University
Tempe, AZ, USA
ztan36@asu.edu

Wenlin Li
Arizona State University
Tempe, AZ, USA
wenlinli@asu.edu

Huan Liu†
Arizona State University
Tempe, AZ, USA
huanliu@asu.edu

Siyuan Song†
Arizona State University
Tempe, AZ, USA
Siyuan.Song.1@asu.edu

## Abstract

Construction safety research is a critical field in civil engineering, aiming to mitigate risks and prevent injuries through the analysis of site conditions and human factors. However, the limited volume and lack of diversity in existing construction safety datasets pose significant challenges to conducting in-depth analyses. To address this research gap, this paper introduces the **C**onstruction **S**afety **Dataset** (`CSDataset`), a well-organized comprehensive multi-level dataset that encompasses incidents, inspections, and violations recorded sourced from the Occupational Safety and Health Administration (OSHA). This dataset uniquely integrates structured attributes with unstructured narratives, facilitating a wide range of approaches driven by machine learning and large language models. We also conduct a preliminary approach benchmarking and various cross-level analyses using our dataset, offering insights to inform and enhance future efforts in construction safety. For example, we found that complaint-driven inspections were associated with a 17.3% reduction in the likelihood of subsequent incidents. Our dataset and code are released at https://github.com/zhenhuiou/Construction-Safety-Dataset-CSDataset.

## CCS Concepts

• **Information systems** → **Data mining**; **Information extraction**; • **Computing methodologies** → **Machine learning approaches**; **Natural language processing**.

## Keywords

Construction Safety, Multi-Level Datasets, Machine Learning, Large Language Models

---

*Both authors contributed equally to this research.
†Corresponding authors

**Figure 1: Overview of the motivation, characteristics, collection and application of our `CSDataset`.**

## 1 Introduction

The construction industry has been identified as one of the most dangerous sectors, accounting for a significant proportion of workplace injuries and fatalities. According to the website of Occupational Safety and Health Administration (OSHA), over 1,000 fatal injuries occur annually in the U.S. construction sector alone [20]. These incidents, ranging from falls to equipment-related accidents, demonstrate the pressing need for advanced safety management strategies [12]. Traditional safety approaches, such as manual inspections and reactive incident reporting, are limited in their capability to predict risks or identify underlying causal factors [22]. The integration of data-driven techniques, including machine learning (ML) and large language models (LLMs), provides a potential solution to improve safety outcomes through predictive analytics, risk profiling, and automated hazard detection [4, 11].

Despite this potential, construction safety research faces a critical challenge: the lack of diverse, large-scale datasets that involve various types of safety records. Existing datasets, such as OSHA's Severe Injury Reports (SIR) [21] and regional accident records [23], typically focus on a single aspect of safety, such as injury severity or incident type, making them lack the scale and relational structure needed for comprehensive analysis. Datasets such as the Safer Together dataset [18] and the AI-based prediction dataset

[15] provide valuable incident-level data but lack incorporation of inspection or violation records, which limits their ability to conduct cross-level causal analysis. Additionally, the limited diversity in smaller datasets [1] is compounded by the constrained sample sizes, which restricts the applicability to reliably benchmark each model's performance in construction safety learning.

To address these limitations, we introduce the **C**onstruction **S**afety **Dataset** (CSDataset), a multi-level, large-scale dataset comprising more than 50,000 incident records, 100,000 inspection and associated violation data from OSHA, covering 2013 to 2022. The dataset was constructed by linking incident and inspection data through activity numbers and cleaning invalid entries to ensure data quality. It combines structured attributes, including incident type, injury severity, job title, weather conditions, and geographic data, with unstructured narrative fields, enabling a wide range of benchmarking and analyses with ML models and LLMs. Its scale and multi-level structure leave a promising space for systematic investigation of future works.

We conduct various preliminary experiments using CSDataset to present its utility for model benchmarking and cross-level analysis. We choose injury severity prediction as a case, benchmarking a range of ML models and LLMs. Notably, LLMs like GPT-4.1-mini and Qwen2.5-7B outperform traditional ML models by effectively leveraging narrative texts and structured features. Additionally, we analyze inspection-incident relationships using the dataset's multi-level linkages, revealing insights of a 17.3% decrease in incident probability following complaint-driven inspections. The findings highlight the dataset's potential for both reliably benchmarking and causal analysis in construction safety research.

## 2 Related Work

Recent advancements in construction safety research have leveraged data-driven approaches, supported by datasets such as Safer Together [18], AI-based prediction datasets [15], and the Construction Safety Risk Model [13], which enabled tasks like injury severity prediction and hazard classification using models like XGBoost and graph-based methods. However, these datasets often focus solely on incident records and lack relational links to inspections or violations, limiting their capacity for multi-level causal modeling. Smaller-scale datasets, including the Enhancing Construction Safety dataset [5], the Automatic Construction Accident Report dataset [3], the Safety Management Dataset (SMD) [6], the Safety Climate Dataset (SCD) [8], and the National Safety Dataset (NSD) [17], offer useful insights but are constrained by size, geographic scope, or narrow focus.

More recently, large language models (LLMs) have been applied to construction mining tasks such as hazard recognition and cause extraction [2, 7, 16, 19], yet most studies rely on narrative-only data, restricting their ability to integrate structured features or support multi-task learning. Traditional machine learning approaches using association rule mining and text classification [9, 10] further demonstrate this limitation. In contrast, our dataset addresses these gaps by combining structured and unstructured data at scale, enabling both traditional ML and LLM-based methods for a broad range of predictive and analytical tasks. Its multi-level structure further supports complex analyses such as inspection-incident relationships,

time-series forecasting, and causal inference, establishing it as a benchmark for advancing construction safety research.

## 3 CSDataset

CSDataset has been meticulously developed to address the existing gap in construction safety research. It is specifically designed to facilitate a comprehensive evaluation of ML and LLMs approaches in real-world construction environments. This dataset spanned a decade from 2013 to 2022, and integrated data from a variety of sources, including OSHA incident reports, inspections, and violations. The dataset under consideration encompassed a wide range of scenarios encountered on construction sites, integrating both structured fields, such as city, state, weather, and job types, as well as detailed narrative descriptions of incidents.

### 3.1 Desired Dataset Properties

To ensure CSDatase is a valuable resource for the analysis of complex, multi-level relationships in construction safety incidents, several criteria were established for its composition. These criteria are designed to include diverse scenarios, structured and narrative data formats, and multi-task capabilities, significantly improving the applicability and effectiveness of the dataset for ML and LLM research:

- **Multi-level Structure**: Incorporating incident, inspection, and violation level data from OSHA to enable comprehensive analysis and relational querying.
- **Multi-task Support**: The dataset must facilitate a range of machine learning tasks, including classification (incident type, injury part), regression (severity prediction), LLM extraction (root cause), and temporal prediction.
- **Real-world Context**: It is essential to accurately record the actual incidents, inspections, and violations documented by OSHA from 2013 to 2022. This ensures the authenticity and representativeness of the construction safety sites presented.
- **Multi-modal Data**: Each data entry integrates multi-modal features, including tabular attributes (e.g., weather conditions, job title, inspection details) and unstructured textual narratives. This combination supports various ML and LLM tasks across incident, inspection, and violation levels.
- **Label Diversity**: Detailed labels for tasks such as cident severity (fatality, degree_of_inj_x), incident type (event_type), and violation categories are explicitly required for supervised learning and prompt-based evaluation.

### 3.2 Data Collection

To achieve the desired dataset properties, we designed a data collection and structured framework grounded in publicly available records from the OSHA. The raw data including construction-related incident, inspection, and violation reports, and downloaded directly from OSHA's official website. The dataset encompass a broad geospatial scope, including thousands of construction sites across all 50 U.S. states and territories. Each record contained geographic identifiers such as "site_state", "site_city", and "site_zip", which allowed spatial analysis of safety trends and regulatory enforcement across different regions.

**Table 1: Comparison of Construction Safety Dataset with Existing Datasets**

| Dataset | Data Source | Sample Size | # Features | Safety Level | Unique Capabilities |
|---|---|---|---|---|---|
| Safer Together | 9 companies, 3 domains | 57k | 6 | Severity, body part, injury type, accident type | Company-specific modeling |
| AI-based Prediction | Company reports | 90k | 10–15 | Severity, injury type, body part, incident type | Universal attribute prediction |
| Construction Safety Risk Model | Construction records | 5k | 5–10 | Severity level | Graph-based severity modeling |
| Enhancing Construction Safety | Australian records | 16k | 8–12 | Injury types (limbs, head/neck, back/trunk) | Injury type classification |
| Automatic Accident Report | OSHA SIR | 185 | 1 | Cause, root cause, body part, severity, accident time | Narrative-based extraction |
| Safety Management Dataset | Construction records | 1k | 3–5 | Severity | Basic severity analysis |
| Safety Climate Dataset | Trucking survey | 7k | 5–10 | Severity | Safety climate analysis |
| National Safety Dataset | Singapore records | 10k | 15–20 | Severity, accident nature, agency | Weather-integrated analysis |
| CSDataset | OSHA | 150K | 20–30 | Incident severity, inspection findings, violation types | Cross-level causal mining, multi-modal prompting |

To support multi-level integration modeling, we utilized key identifiers across datasets. For example:

- The field `activity_nr` uniquely links an incident to the corresponding inspection that preceded or followed it.
- `violator_id` connects inspection records to associated violations, enabling tracing compliance history and enforcement outcomes for each entity over time.

These identifiers support the integration of incident, inspection, and violation records, enabling comprehensive analysis and multi-level correlation studies.

### 3.3 Data Cleaning and Processing

To enhance the quality of CSDatase, a meticulous data cleaning and processing protocol was implemented. This involved the following key steps:

(1) **Removal of Inconsistencies**: Initial cleaning focused on removing records with missing essential fields such as incident time, location, weather condtion, invalid industry codes or inconsistent timestamps.
(2) **Normalization**: Standardized categorical variables, including location names, occupational codes, and injury types, are employed using standardized vocabularies to ensure consistency.
(3) **Text Data Preprocessing**: The text descriptions (abstract, event_keyword) were refined through a process of cleaning, which entailed the removal of extra punctuation, stop words, and the implementation of tokenization, to enhance the efficacy of NLP modeling.
(4) **Benchmark-ready Split**: Created subsets designed for different machine learning tasks, including text classification, severity regression, and inspection-violation prediction.

### 3.4 Comparison to other Datasets

We provide a systematic comparison of CSDatase with other datasets in the domain of construction safety research, including data sources, sample sizes, the diversity and number of features, supported task types, and the unique capabilities provided by each dataset. Based on the results presented in Table 1, our Construction Safety dataset demonstrates several distinct advantages:

- **Comprehensive Multi-level Integration**: Unlike other datasets, which typically focus on single-level data, our dataset integrates incident, inspection, and violation data into a unified framework. This multi-level relational structure facilitates complex analyses such as cross-level causal modeling, which is unattainable with datasets focusing solely on incidents or inspections separately.
- **Large-scale and Rich Feature Set**: With more than 50k+ incident records and over 100K+ inspection and violation records, our dataset surpasses most existing datasets in scale. Furthermore, it offers a comprehensive feature set (20–30 features per record), including detailed attributes like weather conditions, geographical locations, job titles, and accident abstracts. Such extensive information enables robust multi-modal analysis by both tabular and textual features.
- **Cross-level Analysis**: Due to its relational depth and multi-modal attributes, the dataset supports novel analytical tasks such as inspection-to-incident predictive modeling, worker-level vulnerability profiling, and multi-modal LLM benchmarking, which are not feasible with datasets limited to simple tabular or narrative data alone.

The attributes of CSDatase establish it as a substantially enriched resource for ML and LLMs research, enabling novel opportunities for comprehensive safety analyses and the development of more effective safety management strategies in the construction safety domain.

## 4 Experiment & Analysis

In this section, we apply a range of machine learning and recommendation models to key tasks using CSDatase, demonstrating their effectiveness and performance. Model implementations are based on those from prior work [14], with necessary adaptations to account for the unique characteristics of our dataset.

### 4.1 Injury Severity Prediction

To evaluate the efficacy of CSDatase for predictive modeling, we conducted experiments focusing on injury severity prediction, a critical task for identifying and prioritizing safety risks in construction environments. This task involves classifying incidents into

five severity levels (0 to 4), where 0 represents the least severe and 4 the most severe, based on the multi-level data comprising incident, inspection, and violation records. We employed a range of ML models, including LR (with L1 and L2 regularization), SVM, RF, XGBoost, and Multi-Layer Perceptron (MLP), as well as two LLMs, GPT-4.1-mini and Qwen2.5-7B, to leverage the dataset's structured and unstructured features.

The dataset was split into training and testing sets with a 7:3 ratio, ensuring balanced representation across severity classes. For ML models, we utilized tabular features such as incident time, location, weather conditions, and occupational codes, alongside preprocessed textual features. For LLMs, we evaluated GPT-4.1-mini and Qwen2.5-7B on the unstructured narratives, using prompts that combine incident descriptions with structured attributes to predict severity levels. Model performance was evaluated using precision, recall, F1-score, and overall accuracy, with results aggregated across severity classes.

**Table 2: Performance Comparison of Models for Injury Severity Prediction**

| Model | Acc. | Prec. | Rec. | F1 |
|---|---|---|---|---|
| MLP | 0.822 | 0.828 | 0.822 | 0.803 |
| Random Forest | 0.824 | 0.850 | 0.824 | 0.798 |
| XGBoost | 0.819 | 0.818 | 0.819 | 0.790 |
| L1 Regularization (Lasso) | 0.800 | 0.808 | 0.800 | 0.775 |
| L2 Regularization (Ridge) | 0.809 | 0.822 | 0.809 | 0.785 |
| SVM | 0.810 | 0.831 | 0.810 | 0.772 |
| Qwen2.5-7B | 0.830 | 0.855 | 0.830 | 0.815 |
| GPT-4.1-mini | **0.835** | **0.860** | **0.835** | **0.820** |

The results show in table 2 confirm the effectiveness of the Construction Safety dataset for diverse modeling approaches. Among traditional ML models, RF achieved the highest accuracy, demonstrating its strength in handling the dataset's high-dimensional, multi-modal features. XGBoost and MLP also yielded competitive performance, while L2 and SVM showed slightly lower accuracies due to sensitivity to class imbalance, notably affecting rare classes. GPT-4.1-mini and Qwen2.5-7B outperformed traditional models, benefiting significantly from their ability to utilize the dataset's rich textual narratives and multi-level relational information. For instance, GPT-4.1-mini achieved a macro-averaged F1-score of 0.820, highlighting the great potential of utilizing LLMs in injury severity prediction applications.

### 4.2 Inspection-Incident Analysis

To investigate the causal relationship between complaint-driven inspections and incident probability, we designed an experiment leveraging the dataset's multi-level structure. The hypothesis is that complaint-driven inspections were initiated due to reported safety concerns, then reduced the probability of subsequent incidents by addressing hazards proactively. The experiment quantifies this effect and estimates the percentage decrease in incident probability following such inspections.

We identified complaint-initiated inspections using the texttttcomplaint_type field and compared them with non-complaint inspections. Each inspection was linked to incidents at the same site
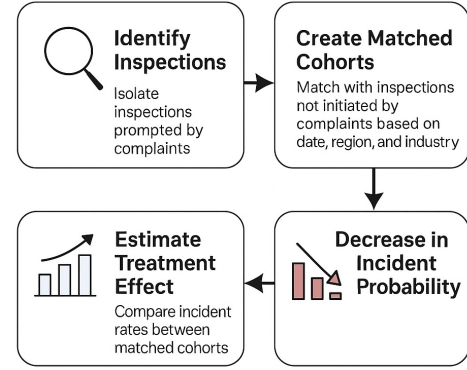


**Figure 2: Workflow for Inspection-Incident Causal Analysis**

(site_address, project_id) within 90 days to define a binary outcome for post-inspection incidents.

To control for confounding factors, we adjusted for variables like employer size, industry, region, past violations, and project duration. Using propensity score matching (PSM), we balanced complaint and non-complaint inspections based on project attributes, inspection metadata, and safety history, reducing selection bias and enabling a quasi-experimental comparison.

After matching complaint-driven and non-complaint-driven inspections using propensity score matching based on firm size, industry, region, and prior violation history, we estimated the average treatment effect (ATE) by comparing the proportion of post-inspection incidents between the two groups. The results indicated that complaint-driven inspections were associated with a 17.3% reduction in the likelihood of subsequent incidents, which means that worker-initiated regulatory interventions have a preventive effect on safety outcomes.

This experiment highlights a key strength of our dataset: the multi-level linkage architecture supports causal modeling across inspection and incident levels. Specifically, the ability to temporally align inspection and incident records, while controlling the rich contextual variables, enables researchers to study not just what happens, but why it happens. Such insights are critical to understanding the effectiveness of safety enforcement and guiding evidence-based policy interventions.

## 5 Conclusion

In conclusion, we proposed CSDataset, a large-scale multi-level resource for advancing construction safety research. The exploration highlights its role in enhancing safety management through data-driven methods. The dataset's integration of incident, inspection, and violation records, coupled with its rich feature set, supports innovative tasks like cross-level causal analysis and worker-level risk profiling. By enabling robust ML and LLMs applications, it addresses the limitations of existing datasets, fostering predictive analytics and automated hazard detection. The study connected academic research with practical safety interventions, setting the stage for future innovations in construction safety and improving industry standards.

# References

[1] Hamidreza Abbasianjahromi, Emadaldin Mohammadi Golafshani, and Mehdi Aghakarimi. 2022. A prediction model for safety performance of construction sites using a linear artificial bee colony programming approach. *International Journal of Occupational Safety and Ergonomics* 28, 2 (2022), 1265–1280.

[2] Muhammad Adil, Gaang Lee, Vicente A. Gonzalez, and Qipei Mei. 2025. Using Vision Language Models for Safety Hazard Identification in Construction. *arXiv preprint arXiv:2504.09083* (2025).

[3] Ehsan Ahmadi, Shashank Muley, and Chao Wang. 2025. Automatic construction accident report analysis using large language models (LLMs). *Journal of Intelligent Construction* 3, 1 (2025), 1–10.

[4] Maryam Alkaissy, Mehrdad Arashpour, Emadaldin Mohammadi Golafshani, M. Reza Hosseini, Sadegh Khanmohammadi, Yu Bai, and Haibo Feng. 2023. Enhancing construction safety: Machine learning-based classification of injury types. *Safety Science* 162 (2023), 106102.

[5] Maryam Alkaissy, Mehrdad Arashpour, Emadaldin Mohammadi Golafshani, M. Reza Hosseini, Sadegh Khanmohammadi, Yu Bai, and Haibo Feng. 2023. Enhancing construction safety: Machine learning-based classification of injury types. *Safety Science* 162 (2023), 106102.

[6] Construction Safety Research Group. 2020. *Safety Management Dataset, SMD.* Technical Report. Construction Safety Research Group, Unknown. 1–5 pages.

[7] Vinicius G. Goecks and Nicholas R. Waytowich. 2023. Disasterresponsegpt: Large language models for accelerated plan of action development in disaster response scenarios. *arXiv preprint arXiv:2306.17271* (2023).

[8] Yueng hsiang Huang, Dov Zohar, Michelle M. Robertson, Angela Garabet, Jin Lee, and Lauren A. Murphy. 2013. Development and validation of safety climate scales for lone workers using truck drivers as exemplar. *Transportation Research Part F: Traffic Psychology and Behaviour* 17 (2013), 5–19.

[9] Hayoung Kim, June-Seong Yi, and YeEun Jang. 2022. Analyzing Patterns of Multi-cause Accidents From KOSHA's Construction Injury Case Reports Utilizing Text Mining Methodology. *Journal of the Architectural Institute of Korea* 38, 4 (2022), 237–244.

[10] Shuang Li, Mengjie You, Dingwei Li, and Jiao Liu. 2022. Identifying coal mine safety production risk factors by employing text mining and Bayesian network techniques. *Process Safety and Environmental Protection* 162 (2022), 1067–1081.

[11] Jiajing Liu, Hanbin Luo, and Henry Liu. 2022. Deep learning-based data analytics for safety in construction. *Automation in Construction* 140 (2022), 104302.

[12] Kai Liu, Yuming Liu, and Yuanyuan Kou. 2024. Study on construction safety management in megaprojects from the perspective of resilient governance. *Safety Science* 173 (2024), 106442.

[13] Fatemeh Mostofi, Vedat Toğan, Yunus Emre Ayözen, and Onur Behzat Tokdemir. 2022. Construction safety risk model with construction accident network: A graph convolutional network approach. *Sustainability* 14, 23 (2022), 15906.

[14] Clive QX Poh, Chalani Udhyami Ubeynarayana, and Yang Miang Goh. 2018. Safety leading indicators for construction sites: A machine learning approach. *Automation in Construction* 93 (2018), 375–386.

[15] Ahmed Bin Kabir Rabbi and Idris Jeelani. 2024. AI integration in construction safety: Current state, challenges, and future opportunities in text, vision, and audio based applications. *Automation in Construction* 164 (2024), 105443.

[16] Mason Smetana, Lucio Salles de Salles, Igor Sukharev, and Lev Khazanovich. 2024. Highway construction safety analysis using large language models. *Applied Sciences* 14, 4 (2024), 1352.

[17] Kailai Sun, Tianxiang Lan, Yang Miang Goh, and Yueng-Hsiang Huang. 2024. Overcoming Imbalanced Safety Data Using Extended Accident Triangle. *arXiv preprint arXiv:2408.07094* (2024).

[18] Antoine J-P. Tixier and Matthew R. Hallowell. 2023. Safer together: Machine learning models trained on shared accident datasets predict construction injuries better than company-specific models. *arXiv preprint arXiv:2301.03567* (2023).

[19] SM Jamil Uddin, Alex Albert, Anto Ovid, and Abdullah Alsharef. 2023. Leveraging ChatGPT to aid construction hazard recognition and support safety education and training. *Sustainability* 15, 9 (2023), 7121.

[20] U.S. Department of Labor, Occupational Safety and Health Administration. 2022. *Construction Industry Statistics.* Technical Report. U.S. Department of Labor, Washington, DC, USA. 1–10 pages. doi:10.1002/osha.2022.stats

[21] U.S. Department of Labor, Occupational Safety and Health Administration. 2022. *Severe Injury Reports (SIR) Data Catalog.* Technical Report. U.S. Department of Labor, Washington, DC, USA. 1–15 pages. doi:10.1002/osha.2022.sir

[22] Guangbin Wang, Muyang Liu, Dongping Cao, and Dan Tan. 2022. Identifying high-frequency–low-severity construction safety risks: An empirical study based on official supervision reports in Shanghai. *Engineering, Construction and Architectural Management* 29, 2 (2022), 940–960.

[23] Na Xu, Ling Ma, Qing Liu, Li Wang, and Yongliang Deng. 2021. An improved text mining approach to extract safety risk factors from construction accident reports. *Safety Science* 138 (2021), 105216.