

DETC2025-169712

**RECALL-MM: A MULTIMODAL DATASET OF CONSUMER PRODUCT RECALLS
FOR RISK ANALYSIS USING COMPUTATIONAL METHODS AND LARGE
LANGUAGE MODELS**

Diana Bolanos

Department of Mechanical Engineering
University of California
Berkeley, California 94720
Email: dbolanos@berkeley.edu

Mohammadmehdi Ataei

Autodesk Research
Toronto, ON M5G 1M1, Canada
Email: mehdi.ataei@autodesk.com

Daniele Grandi

Autodesk Research
San Francisco, CA 94105
Email: danielle.grandi@autodesk.com

Kosa Goucher-Lambert*

Department of Mechanical Engineering
University of California
Berkeley, California 94720
Email: kosa@berkeley.edu

ABSTRACT

Product recalls provide valuable insights into potential risks and hazards within the engineering design process, yet their full potential remains underutilized. In this study, we curate data from the United States Consumer Product Safety Commission (CPSC) recalls database to develop a multimodal dataset, RECALL-MM, that informs data-driven risk assessment using historical information, and augment it using generative methods. Patterns in the dataset highlight specific areas where improved safety measures could have significant impact. We extend our analysis by demonstrating interactive clustering maps that embed all recalls into a shared latent space based on recall descriptions and product names. Leveraging these data-driven tools, we explore three case studies to demonstrate the dataset's utility in identifying product risks and guiding safer design decisions. The first two case studies illustrate how designers can visualize patterns across recalled products and situate new product ideas within the broader recall landscape to proactively anticipate hazards. In the third case study, we extend our approach by employ-

ing a large language model (LLM) to predict potential hazards based solely on product images. This demonstrates the model's ability to leverage visual context to identify risk factors, revealing strong alignment with historical recall data across many hazard categories. However, the analysis also highlights areas where hazard prediction remains challenging, underscoring the importance of risk awareness throughout the design process. Collectively, this work aims to bridge the gap between historical recall data and future product safety, presenting a scalable, data-driven approach to safer engineering design.

1 INTRODUCTION

Risk analysis is a necessary step in the product development process. Engineers and designers are encouraged to predict potential hazards using traditional six-sigma approaches to assess potential failure modes [1]. Nonetheless, consumer products are often recalled due to design and manufacturing related hazards, posing a risk of injury and sometimes death [2]. As such, we see an opportunity to learn from recalled products to observe *what*

* Address all correspondence to this author.

products fail, and *how* the failures occur, ultimately providing engineers and designers historical information of existing failure modes. This study leverages the United States Consumer Product Safety Commission (CPSC) recalls database to serve as a benchmark for novel computational risk prediction approaches presented herein. We curate a dataset of 6,874 recalls spanning dates between the years 2000 and 2024, augmenting the retrieved database information with new descriptors created using a large language model (LLM). Notably, these recalls account for over 546 million individual product SKUs reported as affected over the past two decades, underscoring the vast scale and real-world impact of product safety failure. An example of preprocessed recall entry can be found in Appendix A Table 3. We highlight the use of the dataset and present how it could support risk identification in the design process.

Our contributions include: (1) the development of RECALL-MM, a curated, multimodal dataset of recalled consumer products, augmented through LLM-generated classifications and visual descriptors, (2) the demonstration of computational methods for embedding and visualizing recall data to uncover patterns in product failures, supported by two case studies illustrating how these methods can aid risk identification, and (3) the application of an LLM to predict potential product hazards based solely on visual descriptions, highlighting both the strengths and limitations of automated hazard assessment. Collectively, this work aims to improve product safety, anticipate design failures, and support data-driven decision-making in engineering design.

To support further development and evaluation of our dataset, we make the RECALL-MM dataset and accompanying experimental code publicly available on GitHub¹.

2 RELATED WORK

2.1 Design Datasets

Over the last decade, several large design datasets have been curated and released to support data-driven design efforts for product design and other design tasks. The classes of objects collected, sample size, and modality of the data are the main differentiators between datasets in this field.

Shapenet [3], the ABC Dataset [4], DeepCAD [5], and the Fusion 360 Gallery dataset [6] are among the largest datasets that contain geometry data and class labels for individual parts and whole assemblies. The datasets have been widely used for design automation and geometry generation tasks, and have also supported other work around ancillary design tasks such as materials selection [7]. Other smaller datasets have also been curated around more specific classes of objects, such as car bodies, mechanical components, and bicycles [8–10]. While these datasets provide valuable structured information on product fea-

tures, making them a useful starting point for analysis, they also come with limitations, including incorrect or missing semantic information, limited number of object classes, and lack of design context, intent, and criteria, which limits their utility for comprehensive risk assessment.

Other design datasets have focused on different modalities, such as hand-drawn sketches [11, 12], or textual descriptions of designs [13–16]. Although these datasets focus more on design rationale, describing the features and aesthetics of the design solutions, they do not explicitly consider design feasibility, and some are limited by the number of object classes and data quality.

A few multimodal design datasets have been published, built around graphic design [17], design requirement documents [18], and descriptions of design changes [19]. These datasets combine textual descriptions, geometry, images, and design requirements in various combinations to support design tasks such as editing 3D geometry, interpreting design requirements, and style classification.

The multimodal dataset collected in this work differs from prior datasets as it provides textual descriptions of designs, images, and recall information related to a failure mode of the product over a wide range of product classes.

2.2 Data-Driven Risk Analysis

Product recalls have been used to investigate trends in consumer product safety, with a prevailing focus on children's toys [20–23]. These studies similarly leverage the CPSC database, along with the global recalls dataset from the Organisation for Economic Co-operation and Development (OECD) [24]. Wai and Utama [23] present a series of machine learning approaches for predicting a binary classification of children's toy safety, showcasing the potential for data-driven hazard prediction. Our study expands on this work by investigating trends across multiple product categories, moving beyond a single domain focus.

The automotive industry has also seen a shift towards data-driven design for predictive maintenance and hazard prevention. Yorulmus et al. [25] present machine learning approaches for predicting brake defects from vehicles passing quality assurance checkpoints, yet exhibiting high rates of customer complaints. Similarly, [26] demonstrated a multi-label transfer learning approach by implementing a series of pretrained Convolutional Neural Networks (CNNs) to predict the binary failure status of multiple engine components. Both studies rely on failure data to improve future design of automotive components, demonstrating the effectiveness of leveraging data for failure mitigation and design improvement.

2.3 Hazard Identification in Design

The field of risk analysis has been extensively studied and continues to evolve within various organizational and engineering domains. A fundamental objective of risk analysis in en-

¹<https://github.com/dianabolanos/RECALL-MM>

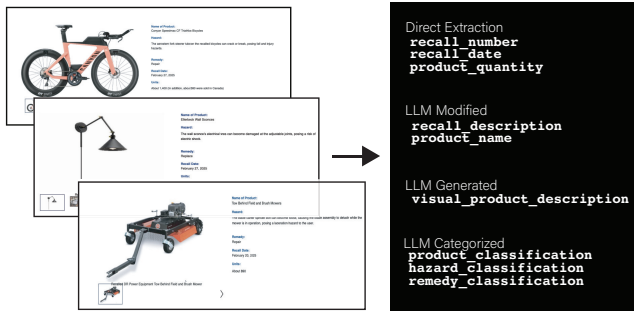


FIGURE 1: Process overview of translating database information into nine distinct data fields.

engineering design is to proactively identify and mitigate hazards before they manifest as failures or safety incidents. Failure Mode and Effects Analysis (FMEA) is among the most widely adopted methodologies employed by organizations to structure systematic risk assessments, prioritize potential hazards, and implement preventive actions [27–29]. First developed in the 1960s by the aerospace industry [30], FMEAs now serve as an industry standard tool across various applications, including automotive design, aerospace engineering, and product development, demonstrating its versatility in supporting product reliability and safety. More recent advancements in FMEA methodologies incorporate computational approaches, such as fuzzy logic, machine learning, and integrated decision-making frameworks, enhancing traditional FMEA practices by addressing uncertainties and subjectivity inherent in risk scoring [31]. These improvements continue to reinforce the importance of hazard identification and risk mitigation in complex engineering designs. We further posit that reviewing historical data can expand the results of risk analysis activities.

3 METHODS

We focus this section on detailing the steps used to clean and augment the CPSC database into a multimodal dataset used throughout the study. Then, we describe the computational methods used to embed this dataset into a vectorized representation, allowing for deeper analysis and visualization. Finally, we introduce the methodology for leveraging an LLM to predict hazards based on visual product information.

3.1 Dataset

The dataset used in this analysis is a curated subset of the US CPSC recalls database [32]. To ensure consistency and feasibility for analysis, recalls were filtered based on API accessibility, the presence of product images, and the timeframe of 2000 to 2024. This filtering process yields a dataset comprising 6,874

recall entries, which accounts for 92.4% of total recalls between the specified timeframe.

Each entry in the dataset includes essential recall attributes such as hazard classifications, product categories, remedy types, historical recall dates, and an associated product image. The raw CPSC data lacked labeled classifications for each entry. To address this, we leveraged generative models, specifically GPT-4o [33], to enrich and structure the data. Figure 1 illustrates which fields were directly extracted from the CPSC database and which were augmented using an LLM.

The **LLM Modified** fields—*recall_description* and *product_name*—were refined using GPT-4o to remove brand references and retain generic descriptions. Additionally, we label the *visual_product_description* field as **LLM Generated**, as it was generated entirely by prompting GPT-4o to describe the product based on its associated image. For classification tasks, GPT-4o assigned each recall to a product category, hazard category, and remedy type, selecting from predefined lists. We label this as **LLM Categorized**. Hazard and remedy categories were aligned with CPSC’s own labels, with definitions provided in Appendix B (Tables 4 and 5). Since the CPSC does not offer standardized product categories, we developed an 11-category scheme based on domain understanding, ensuring it broadly covers the landscape of recalled products while maintaining consistency with OECD terminology [24].

During data cleaning, all records were standardized to a predefined schema, ensuring consistent representation of attributes such as product description, hazard type, and remedy details. GPT-4o outputs were validated against this schema, with type and value constraints applied to ensure reliability. Each record retains its original recall ID, preserving traceability to the source data and supporting future reference or verification.

3.2 Recall Space Exploration and Visualization

3.2.1 Embedding Product descriptors, specifically *recall_description* and *product_name*, were embedded into a numerical latent space using the all-MiniLM-L6-v2 model from Sentence-BERT [34]. Sentence-BERT is a pre-trained model that generates fixed-length dense vector representations optimized for capturing semantic similarity between text inputs. We selected the all-MiniLM-L6-v2 model as it offers an effective balance between model size, computational efficiency, and embedding quality, making it well-suited for large-scale analyses without sacrificing performance.

3.2.2 Dimensionality Reduction for Visualization

To visualize relationships among products and recall reasons, we applied dimensionality reduction techniques. Specifically, we employed the t-distributed Stochastic Neighbor Embedding (t-SNE) algorithm, chosen for its strength in preserving local structure and effectively capturing complex, non-linear relationships

in high-dimensional data [35]. Compared to linear methods such as Principal Component Analysis (PCA) [36], which primarily maintain global variance, t-SNE excels at revealing dense clusters and neighborhood groupings, which are paramount features for identifying semantically similar products and localized recall patterns. Product embeddings obtained via Sentence-BERT were projected from their original vector space into two and three-dimensional coordinates. The resulting visual maps (Fig.5 and Fig.6) enable exploration of recall clusters, offering insight into product similarities and hazard trends. While other methods like UMAP [37] could also be considered, we prioritized t-SNE for its well-established use in exploratory visualizations where fine-grained local structure is of primary interest.

3.3 LLM-Based Hazard Prediction

In addition to computational exploration of the recall data, we evaluate the feasibility of using LLMs to predict potential hazards directly from product images. Specifically, we focus on the *visual_product_description* field, which contains a textual description of each product image generated by GPT-4o.

To perform hazard prediction, we prompt an LLM to analyze the textual description of the image and output all applicable hazard classifications. The model selects hazard labels from a predefined set of ten hazard categories (Appendix B, Table 4), ensuring consistency with existing CPSC classifications. The full prompt used for LLM prediction is provided:

```
You are a product safety expert.
Identify all potential hazards
for the given product. Provide
output in valid JSON format only,
structured as:

{
  'product': {product_description},
  'predicted_hazards': [List of all
applicable hazards]
}

The predicted_hazards field must
only contain hazards from this set:
{all_hazards}. Do not leave any
predicted_hazards fields empty. If
multiple hazards apply, include all
relevant ones. No explanations -
return JSON only.
```

Outputs are returned in a strict JSON format, listing the product description and the predicted hazards. The prompt enforces that the model must not leave any hazard fields empty, encouraging comprehensive identification of potential risks.

Evaluation Metric To quantify the model's performance, we introduce a Relaxed Accuracy (RA) metric:

$$\delta_i = \begin{cases} 1, & \text{if } g_i \in P_i \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

$$\text{Relaxed Accuracy (RA)} = \frac{1}{N} \sum_{i=1}^N \delta_i \quad (2)$$

where

N = Total number of products per hazard class.

g_i = Single ground truth hazard classification for product i .

P_i = Set of predicted hazard classifications for product i .

δ_i = Indicator function that equals 1 if the ground truth hazard g_i is present in the predicted set P_i , and 0 otherwise.

The RA metric accounts for cases where the LLM predicts multiple hazards, but the recall dataset provides only one ground truth hazard per entry. The recall dataset only labels a single hazard per product, even though products may exhibit multiple concurrent hazards. Thus, RA prioritizes capturing all realistic hazards that could apply to a product, rather than penalizing extra predictions, acknowledging that the recall labels may not list every possible hazard.

4 RESULTS AND DISCUSSION

To evaluate the reliability of the curated and augmented dataset, we begin by validating the LLM-generated categorizations against human-annotated ground truths, ensuring consistency across product, hazard, and remedy classifications. Following this, we analyze aggregate trends within the dataset, examining prevalent hazards, product categories, and remedy actions over time. Building on these observations, we then present three case studies to illustrate different approaches for leveraging the dataset in risk identification and design decision-making. The first two case studies employ computational methods, embedding recall data into a shared latent space to explore product relationships and potential risks. The third case study investigates the feasibility of using LLMs to predict potential hazards based solely on visual context.

4.1 Human Evaluation of Dataset Categorizations

To validate the reliability of the LLM categorizations, we compare them against ground truth labels derived from three independent human annotators. Each annotator performed 100

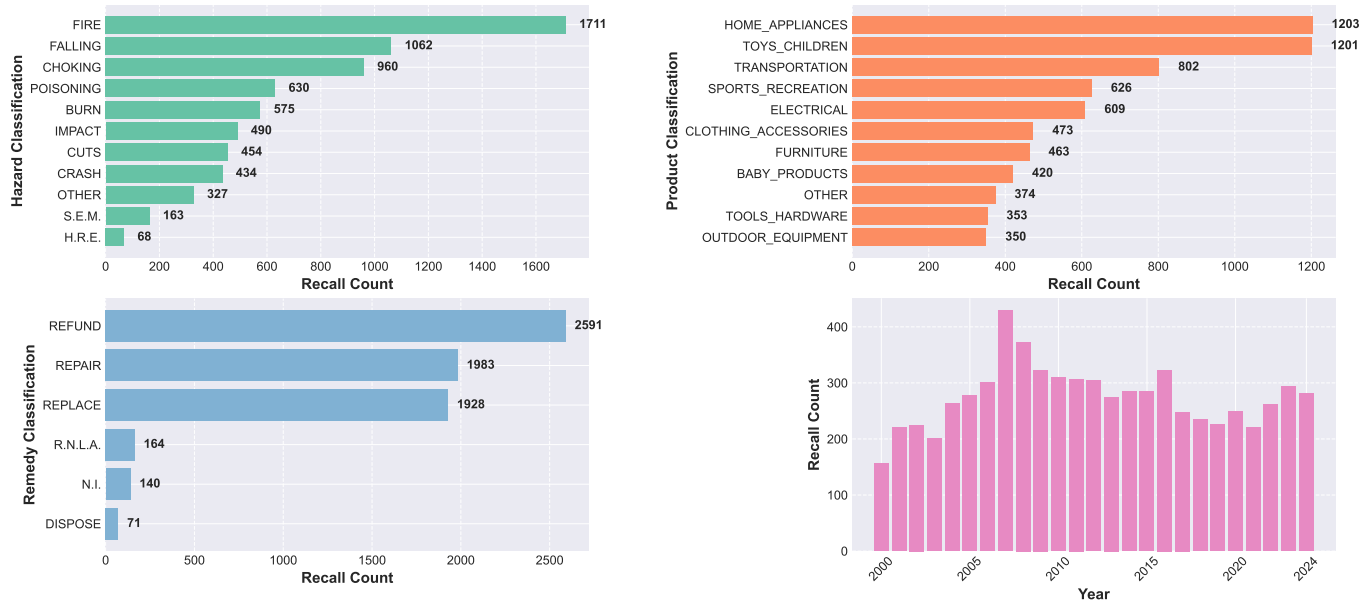


FIGURE 2: Data metrics from 6,874 recalls spanning 2000 - 2024 recall dates.

classifications for each task, amounting to 900 annotations. To establish a ground truth, we employed majority voting across the three annotators' labels for each classification task (product, hazard, and remedies). We assessed inter-rater reliability by evaluating Fleiss's Kappa, which yielded coefficients of 0.71 for product classification, 0.80 for hazard classification, and 0.85 for remedies classification. These scores indicate substantial to almost perfect agreement, based on standard interpretation thresholds. Given this high level of consistency, majority voting was deemed appropriate to consolidate the annotations. Items where no majority agreement was reached (5 for product, 4 for hazard, and 0 for remedy) were excluded from further analysis to maintain the integrity of the ground truth labels.

With the ground truths labels established, we now compare against the LLM categorizations. Using Cohen's Kappa, we observed almost perfect agreement across all three classification tasks, with coefficients of 0.82 for product classification, 0.91 for hazard classification, and 0.90 for remedies classification. These strong agreement levels indicate that the LLM's predictions align closely with human judgment, achieving a level of consistency comparable to expert annotators. Given these results, we are confident in the robustness and accuracy of the LLM-generated outputs and proceed to use them for subsequent analyses in this paper.

4.2 Analysis of Recall Classifications

To analyze patterns in product recalls, we first aggregated the dataset across four dimensions: hazard classifications, product categories, remedy classifications, and recall year. These aggregations are visualized in Figure 2, allowing us to identify prevalent hazards, frequently recalled product types, common industry remedies, and temporal trends in recall activity.

Additionally, to explore relationships between product categories and associated hazards, we generated a hazard-product co-occurrence heatmap (Figure 3). This visualization highlights where certain hazards are disproportionately concentrated within specific product types, offering insight into recurring failure modes and industry-specific safety concerns.

Examining **hazard classification**, *fire*, *falling* and *choking* emerged as the most prevalent hazards, suggesting a need for improved foresight in anticipating these failure modes. Focusing on *fire*, we also see from Fig. 3 that the strongest correlations come from *electrical* and *home_appliances*. This aligns with existing literature indicating heightened risks of household fires due to electrical failures in sockets, plugs, and wiring, as opposed to householder carelessness [38].

Within **product classification**, the high frequency of recalls in *home_appliances* and *toys_children* indicates particular vulnerability to hazards faced within the average US household, indicating the pressing need for safer design of products intended

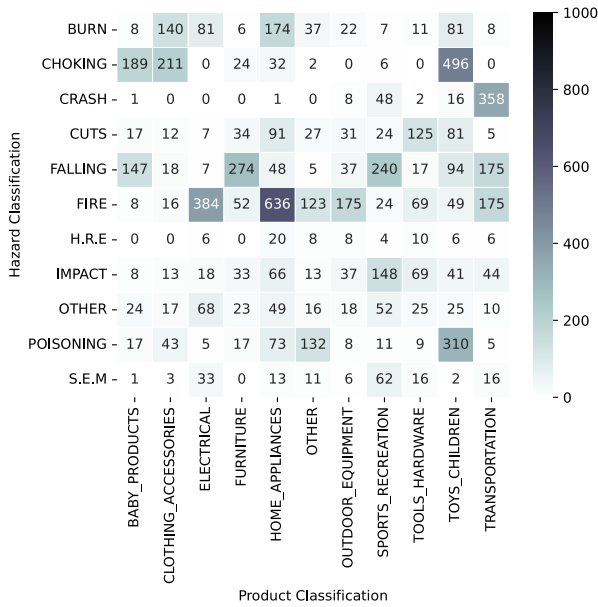


FIGURE 3: Correlation matrix of ground truth product and hazard classifications.

for vulnerable or high-use demographics. A study conducted by Anwar [39] also relied on the CPSC database to examine the harms resulting from high recalls in the toy industry. This study found that while most toys were manufactured in China, a vast quantity of toys were designed in the US, leading to harms related to choking and lead poisoning as primary concerns. This analysis aligns with our findings, emphasizing the importance of stronger safety precautions when designing consumer products. Interestingly, the heatmap also shows lower recall frequencies for categories such as *tools hardware* and *outdoor equipment*, suggesting a possibility of heightened risk awareness and more conservative design practices within these industries.

The prevalence of *refund* as a remedial action indicates a preference towards immediate consumer safety, likely chosen when repairs or replacements are insufficient for risk mitigation. The substantial use of *repair* and *replace* remedies further suggests a widespread industry practice of addressing safety concerns through corrective product interventions rather than solely financial compensation. Kubler et al. [40] conducted a study to observe the effect of product recalls on brand loyalty. Results found that consumers valued transparency and convenient handling of the recalled product.

The temporal analysis revealed an apparent peak in recall incidents around 2005, followed by a general downward trend with fluctuations thereafter. This trend may reflect enhanced regulatory interventions, evolving industry standards, or changes in manufacturing practices and quality control measures. Notably,

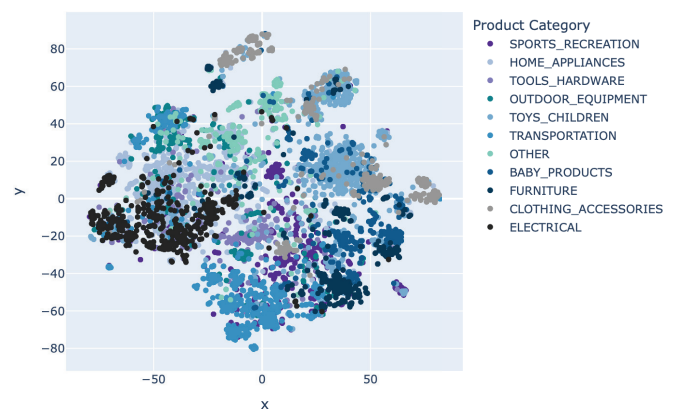


FIGURE 4: Embedding space of recall descriptions labeled by product categories.

the recent uptick observed post-2020 signals the potential effects of disruptions in supply chains due to global events.

Taken together, these results emphasize key areas for improved product safety. Recalls in categories such as *home appliances*, *electrical*, and *toys children* underscore the need for more proactive hazard anticipation during product design. Furthermore, the correlation patterns revealed by the heatmap suggest potential value in cross-domain learning: designers may benefit from examining hazard trends in adjacent product sectors to better anticipate potential risks.

4.2.1 Case Study 1: 2D Latent Space of Recall Descriptions In this case study, we examine the embedding space of the *recall_description* field to explore how textual recall data can reveal underlying patterns beyond predefined hazard classifications. We specifically chose to embed *recall_description* to augment and challenge existing groupings, offering an opportunity to uncover nuanced relationships within the dataset. Figure 4 presents a latent space representation of all 6,874 *recall_description* texts, revealing natural clusters that broadly correspond to different product categories.

To illustrate how product domains influence the structure of this space, we highlight two specific examples in Figure 5. In Figure 5a, we focus on the product categories *electrical* (black) and *clothing accessories* (gray). Here, we observe minimal overlap in the embedding space, reflecting distinct recall descriptions and associated hazards. For example, recalls within the *electrical* category predominantly reference fire hazards, resulting in a more cohesive clustering pattern. In contrast, the *clothing accessories* category exhibits multiple distinct clusters, reflecting a wider variety of recall reasons, such as choking hazards due to detachable components. This divergence underscores how

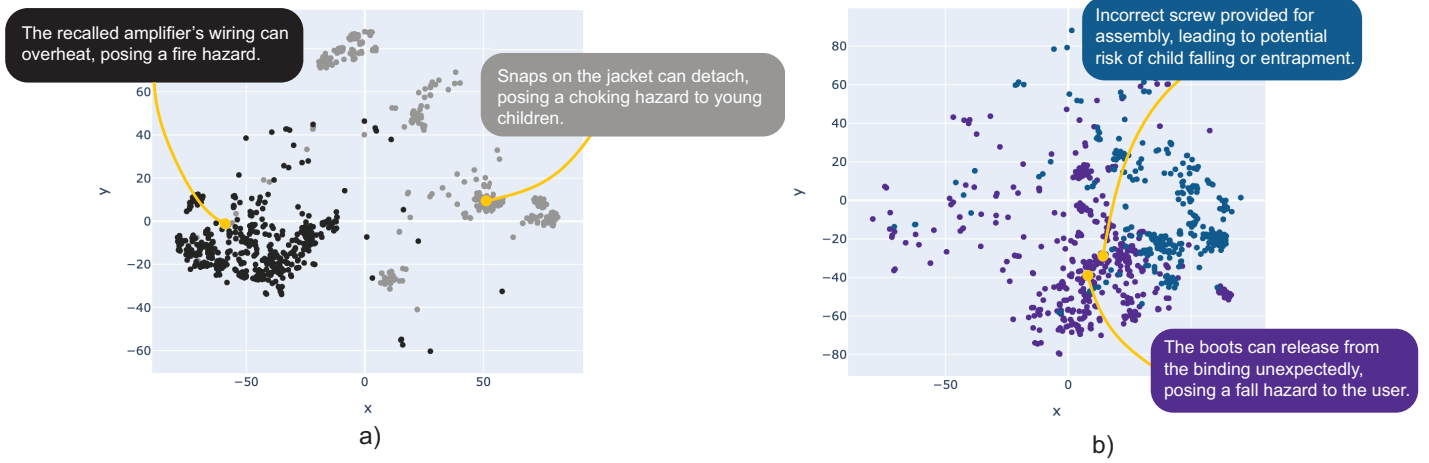


FIGURE 5: a) Embedded recall descriptions of *electrical* (black) and *clothing_accessories* (gray). b) Embedded recall descriptions of *baby_products* (blue) and *sports_recreation* (purple). Descriptions are denoted for each example, showcasing distant (a) and near (b) recall descriptions across different product categories.

certain domains exhibit unique, domain-specific risks, while others exhibit larger diversity in risk possibilities.

Figure 5b highlights a case where recall descriptions from different domains show considerable overlap. Specifically, the categories *baby_products* and *sports_recreation* share similarities in recall descriptions, often referencing issues related to unsecured assembly constraints that pose risks of falling or en-

TABLE 1: Recall description diversity of each product category as measured by Convex Hull area in 2D scaled embedding space.

Category	Normalized Area
TOYS_CHILDREN	11.848
SPORTS_RECREATION	11.009
TRANSPORTATION	10.486
HOME_APPLIANCES	10.169
TOOLS_HARDWARE	9.776
OUTDOOR_EQUIPMENT	9.485
BABY_PRODUCTS	9.219
FURNITURE	9.111
OTHER	8.872
CLOTHING_ACCESSORIES	8.712
ELECTRICAL	6.688

trapment. Despite the differing nature of these product types, the commonality in risk profiles suggests valuable cross-domain learning opportunities. Designers working in either space could benefit from studying hazards in the other, enabling a more comprehensive approach to safety.

To quantitatively contextualize these findings, we calculate the spread of recall descriptions for each product category using Convex Hull analysis. Table 1 reports the normalized areas for each category, offering a metric for the diversity of recall descriptions. This approach, informed by prior work on diversity metrics in embedding spaces [41], enables a comparative assessment of risk variability across categories. Notably, the *toys_children* category exhibits the largest normalized area, indicating a wide range of distinct hazards associated with children's products. This reinforces the need for cautious safety considerations in the design of children's toys.

4.2.2 Case Study 2: 3D Latent Space of Product Name We demonstrate the potential to move beyond passive data exploration by developing an interactive visualization tool that situates new product ideas within the context of historical recall data (see Fig. 6). This visualization embeds all product descriptors—specifically *product_name*—into a shared three-dimensional latent space using t-SNE. When a designer inputs a new product concept, the system embeds the input text into the same space and projects it alongside past recalled products. We chose to embed *product_name* as it typically conveys precise, yet high-level semantic information, making it particularly suitable for early-stage ideation.

For instance, an engineering team might propose a product

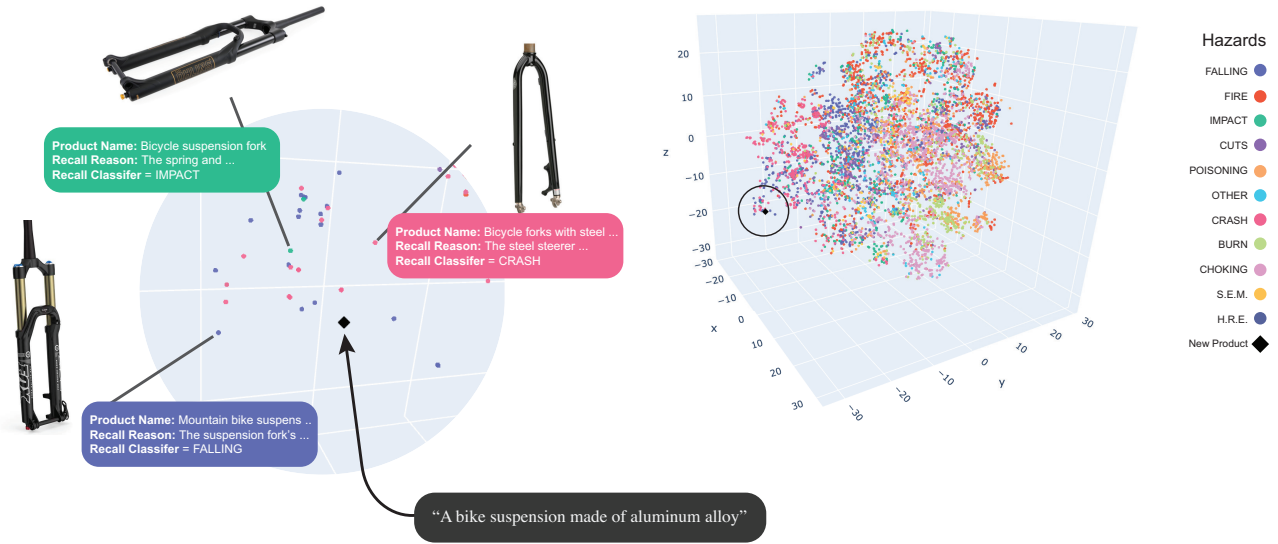


FIGURE 6: 3D embedding space of product name colored by hazard class. An example is shown identifying similarities between a new product and existing products with recall information provided.

idea described simply as “a bike suspension made of aluminum alloy” without fully developed specifications. This tool allows them to position that idea within the landscape of similar historical recalls, providing insight into neighboring products and their associated hazard classifications. By visualizing these relationships interactively, designers can proactively identify potential risks, explore relevant precedents, and make more informed design decisions. Ultimately, this integration of recall data into the early design process supports safer, more responsible product development.

4.2.3 Case Study 3: LLM Hazard Prediction In the third case study, we evaluate the ability of LLMs to predict potential hazards based solely on a product’s text description of the image. Using the approach detailed in Section 3.3, the LLM analyzes the *visual_product_description* field and outputs hazard classifications drawn from the predefined list of categories.

The aggregated frequencies of the LLM-predicted hazard classes are shown in Fig. 7. To assess predictive performance, we compute the RA metric (Eq. (2)) across all hazard classes. The per-class RA scores are summarized in Table 2, demonstrating strong overall predictive performance with an overall RA of 0.73. The results indicate strong predictive capabilities within several hazard categories, with particularly high RA scores observed for *choking* (0.93) and *crash* (0.91) hazards. Conversely, the poisoning hazard class yields the lowest RA score of 0.32.

The relatively low RA scores for certain classes suggest that the model exhibits selectivity in its predictions and is not simply assigning all possible hazards to every product.

TABLE 2: Accuracy per class and overall relaxed accuracy.

Hazard Classification	Relaxed Accuracy
BURN	0.59
CHOKING	0.91
CRASH	0.93
CUTS	0.65
FALLING	0.85
FIRE	0.74
H.R.E.	0.74
IMPACT	0.87
OTHER	0.46
POISONING	0.32
S.E.M.	0.49
Overall Relaxed Accuracy	0.73

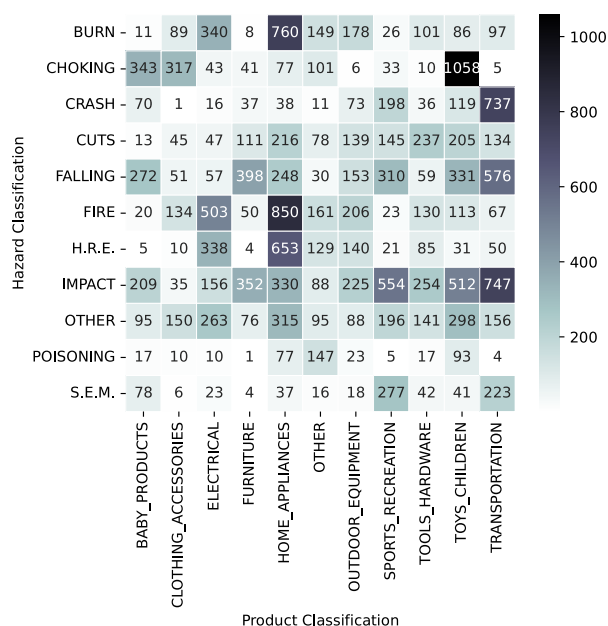


FIGURE 7: Heatmap showing correlations between LLM predicted product and hazard recalls.

Further examination of Figure 7 reveals that the model rarely predicts poisoning hazards for products such as children’s toys. However, cross-referencing with actual recall data (Fig. 3) shows a substantial number of poisoning-related recalls, particularly within the children’s toys category. This discrepancy highlights a critical limitation: certain hazards, like poisoning, may not be visually apparent and thus are underrepresented in LLM predictions. The model’s difficulty in predicting non-visible hazards mirrors how consumers often rely on visual inspection to assess product safety. This highlights a critical need for both transparent hazard communication and proactive design strategies that address risks not readily apparent through appearance alone, particularly in preventing latent hazards such as poisoning.

5 DATASET USE CASES AND FUTURE APPLICATIONS FOR ENGINEERING DESIGN

There are several promising avenues to extend and apply this work. First, augmenting the CPSC dataset with globally recalled product data, such as the OECD Global Recalls Portal [24], would allow for a more comprehensive cross-national analysis. This could reveal broader patterns in product failures and facilitate comparative studies across regulatory environments and cultural contexts.

Additionally, while the methods presented demonstrate feasibility in organizing and classifying recall data, validating their

effectiveness in real-world design and safety processes remains an open opportunity. Future studies could engage with industry practitioners to assess how historical recall data—structured and augmented as shown here—can be integrated into existing risk analysis workflows. Specifically, measuring where practitioners derive value would provide actionable insights.

The multimodal nature of the dataset could also support additional work with visual language models (VLMs). VLMs could be used to predict a product’s hazard and recall risk, similar to Case Study 3 but with additional visual cues from the product’s image. Future studies could determine whether visual or textual cues provide stronger signals for this task. Since directly feeding the image to the VLM for hazard prediction may yield different results, future work could compare these two pipelines systematically, assessing whether the intermediate textual representation improves hazard specificity and alignment with human judgment.

In addition to analyzing past recalls, this dataset and classification pipeline could also be used proactively to improve future product design. For instance, integrating recall-informed hazard classifications into early-stage product requirement generation may improve safety considerations from the outset. This could involve coupling the dataset with LLM generated user requirements [42] to ensure risk factors identified in past recalls are embedded into design requirements.

The dataset can further support developing detailed user personas based on specific recall incidents and hazard types (integrating into works such as [42]). By examining product recalls through the lens of potential user interactions, designers could construct user personas representative of individuals most likely to encounter or exacerbate certain hazards. For example, examining products recalled due to choking hazards might inform the creation of personas representing families with young children or elderly individuals with limited mobility, allowing designers to conduct roleplay analysis and proactively consider how different user behaviors and demographics might interact with products to induce hazardous situations.

Although it was not used in this work, the *remedy_classification* metadata included in the dataset could be used to train a model useful to design practitioners in determining an appropriate solution once a recall-level hazard has been found. Further, future work could explore embedding additional multimodal fields, such as the *visual_product_description*, into a unified vector space. Doing so would support more detailed similarity analyses, helping designers quickly identify potential risks based on visual and textual product features. Investigating the downstream implications of these embeddings – such as their application in automated hazard detection or AI-augmented design tools – presents another valuable research direction.

Ultimately, by refining these methods and integrating them into decision-making pipelines, this research could contribute to addressing future challenges in product safety, thereby encourag-

ing a practice of proactive, data-driven risk mitigation in product design and development.

6 CONCLUSION

This work demonstrates the feasibility and value of leveraging historical product recall data to identify potential hazards in consumer products. By analyzing recall records, we aim to provide designers and engineers with actionable insights into common failure modes and safety risks, ultimately informing safer and more robust product development. Specifically, we advocate for the integration of publicly available datasets into the early stages of the design process, where risk identification is often most critical yet under-informed.

One of the key challenges in early-stage design is anticipating latent hazards that may not be immediately apparent. Through three distinct case studies, we illustrate the utility of computational and LLM-driven methods for interacting with the dataset. These case studies highlight different modalities of engagement: analyzing textual recall descriptions, embedding product names for similarity assessment, and predicting potential hazards from images. These studies demonstrate approaches for understanding not only *what* products fail, but *how* those failures manifest.

By integrating historical recall data into the design process, we present a scalable and data-driven approach to improve product safety, anticipate failure modes, and support risk-informed decision-making. This research lays the groundwork for future efforts aimed at embedding recall-informed analyses into design workflows, ultimately fostering proactive and data-supported risk mitigation in engineering design.

ACKNOWLEDGMENT

Authors DB and KGL would like to acknowledge financial support from the Society of Hellman Fellows and Autodesk Research. Any opinions, findings, and conclusions or recommendations expressed in this article are those of the authors and do not necessarily reflect the views of the aforementioned sponsors.

REFERENCES

- [1] Ulrich, K. T., and Eppinger, S. D., 2016. *Product design and development*. McGraw-hill.
- [2] Hora, M., Bapuji, H., and Roth, A. V., 2011. “Safety hazard and time to recall: The role of recall strategy, product defect type, and supply chain player in the us toy industry”. *Journal of Operations Management*, **29**(7-8), pp. 766–777.
- [3] Chang, A. X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., and Yu, F., 2015. ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago.
- [4] Koch, S., Matveev, A., Jiang, Z., Williams, F., Artemov, A., Burnaev, E., Alexa, M., Zorin, D., and Panozzo, D., 2019. “Abc: A big cad model dataset for geometric deep learning”. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [5] Wu, R., Xiao, C., and Zheng, C., 2021. “Deepcad: A deep generative network for computer-aided design models”. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6772–6782.
- [6] Willis, K. D., Jayaraman, P. K., Chu, H., Tian, Y., Li, Y., Grandi, D., Sanghi, A., Tran, L., Lambourne, J. G., Solar-Lezama, A., et al., 2022. “Joinable: Learning bottom-up assembly of parametric cad joints”. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15849–15860.
- [7] Bian, S., Grandi, D., Liu, T., Jayaraman, P. K., Willis, K., Sadler, E., Borijin, B., Lu, T., Otis, R., Ho, N., et al., 2024. “Hg-cad: hierarchical graph learning for material prediction and recommendation in computer-aided design”. *Journal of Computing and Information Science in Engineering*, **24**(1), p. 011007.
- [8] Kim, S., Chi, H.-g., Hu, X., Huang, Q., and Ramani, K., 2020. “A large-scale annotated mechanical components benchmark for classification and retrieval tasks with deep neural networks”. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, Springer, pp. 175–191.
- [9] Regenwetter, L., Curry, B., and Ahmed, F., 2022. “Biked: A dataset for computational bicycle design with machine learning benchmarks”. *Journal of Mechanical Design*, **144**(3), p. 031706.
- [10] Elrefaie, M., Morar, F., Dai, A., and Ahmed, F., 2024. “Drivaernet++: A large-scale multimodal car dataset with computational fluid dynamics simulations and deep learning benchmarks”. In *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, eds., Vol. 37, Curran Associates, Inc., pp. 499–536.
- [11] Toh, C. A., and Miller, S. R., 2016. “Creativity in design teams: the influence of personality traits and risk attitudes on creative concept selection”. *Research in Engineering Design*, **27**(1), Jan., p. 73–89.
- [12] Gryaditskaya, Y., Sypsteyn, M., Hoftijzer, J. W., Pont, S., Durand, F., and Bousseau, A., 2019. “Opensketch: a richly-annotated dataset of product design sketches”. *ACM Trans. Graph.*, **38**(6), Nov., pp. 232:1–232:16.
- [13] Zhu, Q., and Luo, J., 2022. “Generative pre-trained transformer for design concept generation: An exploration”. *Proceedings of the Design Society*, **2**, May, p. 1825–1834.


- [14] Meltzer, P., Lambourne, J. G., and Grandi, D., 2023. “What’s in a name? evaluating assembly-part semantic knowledge in language models through user-provided names in computer aided design files”. *Journal of Computing and Information Science in Engineering*, **24**(011002), June.
- [15] Jiang, S., Sarica, S., Song, B., Hu, J., and Luo, J., 2022. “Patent data for engineering design: A critical review and future directions”. *Journal of Computing and Information Science in Engineering*, **22**(060902), Oct.
- [16] Goucher-Lambert, K., and Cagan, J., 2019. “Crowdsourcing inspiration: Using crowd generated inspirational stimuli to support designer ideation”. *Design Studies*, **61**, Mar., p. 1–29.
- [17] Lin, J., Huang, D., Zhao, T., Zhan, D., and Lin, C.-Y., 2024. “Designprobe: A graphic design benchmark for multimodal large language models”. arXiv:2404.14801 [cs].
- [18] Doris, A. C., Grandi, D., Tomich, R., Alam, M. F., Ataei, M., Cheong, H., and Ahmed, F., 2025. “Designqa: A multimodal benchmark for evaluating large language models’ understanding of engineering documentation”. *Journal of Computing and Information Science in Engineering*, **25**(2), p. 021009.
- [19] Achlioptas, P., Huang, I., Sung, M., Tulyakov, S., and Guibas, L., 2023. “ShapeTalk: A language dataset and framework for 3d shape edits and deformations”. In Conference on Computer Vision and Pattern Recognition (CVPR).
- [20] Durrett, J., 2015. “A decade of data: an in-depth look at 2014 and a ten-year retrospective on children’s product recalls”. *Kids In Danger*.
- [21] Niven, C. M., Mathews, B., Harrison, J. E., and Vallmuur, K., 2020. “Hazardous children’s products on the Australian and US market 2011–2017: an empirical analysis of child-related product safety recalls”. *Injury prevention*, **26**(4), pp. 344–350.
- [22] Kirschman, K. B., and Smith, G. A., 2007. “Resale of recalled children’s products online: an examination of the world’s largest yard sale”. *Injury Prevention*, **13**(4), pp. 228–231.
- [23] Wai, H. T., and Uttama, S., 2024. “Effectiveness of classifying unsafe children’s toys using nlp, deep learning and ensemble learning”. In 2024 21st International Joint Conference on Computer Science and Software Engineering (JC-SSE), IEEE, pp. 261–267.
- [24] Organisation for Economic Co-operation and Development (OECD), n.d.. Global recalls portal. <https://globalrecalls.oecd.org/#/>. Accessed: 2025-03-08.
- [25] Yorulmuş, M. H., Bolat, H. B., and Bahadır, Ç., 2022. “Predictive quality defect detection using machine learning algorithms: A case study from automobile industry”. In Intelligent and Fuzzy Techniques for Emerging Conditions and Digital Transformation: Proceedings of the INFUS 2021 Conference, held August 24–26, 2021. Volume 2, Springer, pp. 263–270.
- [26] O’Donnell, J., and Yoon, H.-S., 2024. “Determination of multi-component failure in automotive system using deep learning”. *Journal of Computing and Information Science in Engineering*, **24**(2).
- [27] Stamatis, D. H., 2003. *Failure mode and effect analysis: FMEA from theory to execution*. Quality Press.
- [28] Carlsson, C. S., 2012. *Effective FMEAs: achieving safe, reliable, and economical products and processes using failure mode and effects analysis*. John Wiley & Sons.
- [29] Modarres, M., 2006. *Risk analysis in engineering: techniques, tools, and trends*. CRC press.
- [30] Bowles, J. B., and Peláez, C. E., 1995. “Fuzzy logic prioritization of failures in a system failure mode, effects and criticality analysis”. *Reliability engineering & system safety*, **50**(2), pp. 203–213.
- [31] Liu, H.-C., Liu, L., and Liu, N., 2013. “Risk evaluation in failure mode and effects analysis with extended vikor method under fuzzy environment”. *Expert Systems with Applications*, **40**(5), pp. 1795–1803.
- [32] U.S. Consumer Product Safety Commission, 2025. CPSC Recalls Database. Accessed: 2025-03-04.
- [33] OpenAI, 2024. Gpt-4o. <https://openai.com/gpt-4o>. Accessed: 2025-03-21.
- [34] Reimers, N., and Gurevych, I., 2019. “Sentence-bert: Sentence embeddings using siamese bert-networks”. *arXiv preprint arXiv:1908.10084*.
- [35] Van der Maaten, L., and Hinton, G., 2008. “Visualizing data using t-sne”. *Journal of machine learning research*, **9**(11).
- [36] Abdi, H., and Williams, L. J., 2010. “Principal component analysis”. *Wiley interdisciplinary reviews: computational statistics*, **2**(4), pp. 433–459.
- [37] McInnes, L., Healy, J., and Melville, J., 2018. “Umap: Uniform manifold approximation and projection for dimension reduction”. *arXiv preprint arXiv:1802.03426*.
- [38] Taylor, M. J., Fielding, J., and O’Boyle, J., 2024. “Electrical home fire injuries analysis”. *Fire*, **7**(12), p. 471.
- [39] Anwar, S. T., 2014. “Product recalls and product-harm crises: A case of the changing toy industry”. *Competitiveness Review: An International Business Journal incorporating Journal of Global Competitiveness*, **24**(3), pp. 190–210.
- [40] Kübler, R. V., and Albers, S., 2012. “The impact of product recall communication on brand image, brand attitude and perceived product quality”. *Brand Attitude and Perceived Product Quality (May 16, 2012)*.
- [41] Ma, K., Grandi, D., McComb, C., and Goucher-Lambert, K., 2025. “Do large language models produce diverse

design concepts? a comparative study with human-crowdsourced solutions”. *Journal of Computing and Information Science in Engineering*, **25**(2).

- [42] Ataei, M., Cheong, H., Grandi, D., Wang, Y., Morris, N., and Tessier, A., 2025. “Elicitron: A large language model agent-based simulation framework for design requirements elicitation”. *Journal of Computing and Information Science in Engineering*, **25**(2).

Appendix A

TABLE 3: Example of a recall entry from the curated dataset.

Field	Value
recall_number	14259
recall_date	2014-08-20
recall_description	The length adjustment buckles release unexpectedly, causing the item being stored to fall and injure people nearby.
product_name	Kayak and watersports storage hanger
product_quantity	10,000
remedies	Consumers should stop using the recalled storage hangers and return them to the place of purchase for a full refund or replacement.
visual_product_description	The product consists of a pair of straps made from blue and black fabric, each approximately 1 inch wide and 84 inches long when unbuckled. They feature plastic snap buckles for length adjustment and plastic-coated steel S-hooks for hanging.
product_classification	SPORTS_RECREATION
hazard_classification	FALLING
remedies_classification	REPLACE
product_image	

Appendix B

TABLE 4: Hazard classifications and definitions.

Hazard Classification	Definition
Fire	Use of the product may lead to a fire or the product violates federal fabric flammability regulations.
Burn	Use of the product may lead to experiencing burns.
Heat-Related Explosion (H.R.E.)	The product may explode unintentionally.
Falling	Use of the product may cause an unintentional fall.
Poisoning	Use of the product may lead to poisoning.
Crash	Use of the product may lead to an unintentional crash.
Choking	Use of the product may lead to choking, or the product violates federal toy safety standards, or the product violates federal children clothing standards (drawstrings).
Cuts	Use of the product may lead to unintentional cuts and/or lacerations.
Safety Equipment Malfunction (S.E.M.)	The safety product does not operate as intended and use of the product may lead to injury or death.
Impact	Use of the product may lead to an unintentional impact that may cause injury or death.

TABLE 5: Remedy classifications and definitions.

Remedy Classification	Definition
Refund	A customer may receive a full or partial refund, or gift card for the recalled product.
Repair	The company is offering a repair to the recalled product.
Replace	The company is offering a replacement for the recalled product in the form of a new product or other products of similar value.
Dispose	The product should be thrown out or recycled.
New Instructions (N.I.)	The company will issue new instructions on how the customer can make the recalled product safe.
Remedy No Longer Available (R.N.L.A.)	The recalled product should be thrown out or recycled.