

LLMs as World Models: Data-Driven and Human-Centered Pre-Event Simulation for Disaster Impact Assessment

Anonymous ACL submission

Abstract

Efficient simulation is essential for enhancing proactive preparedness for sudden-onset disasters such as earthquakes. Recent advancements in large language models (LLMs) as world models show promise in simulating complex scenarios. This study examines multiple LLMs to proactively estimate perceived earthquake impacts. Leveraging multimodal datasets including geospatial, socioeconomic, building, and street-level imagery data, our framework generates Modified Mercalli Intensity (MMI) predictions at zip code and county scales. Evaluations on the 2014 Napa and 2019 Ridgecrest earthquakes using USGS “Did You Feel It? (DYFI)” reports demonstrate significant alignment, as evidenced by high correlation of 0.88 and low RMSE of 0.77 as compared to real reports at the zip code level. Techniques such as RAG and ICL can improve simulation performance, while visual inputs notably enhance accuracy compared to structured numerical data alone. These findings show the promise of LLMs in simulating disaster impacts that can help strengthen pre-event planning.

1 Introduction

Natural disasters often disrupt infrastructure, causing significant human and economic losses (Jones et al., 2022). Efficient impact assessment is critical for emergency response and evaluating community resilience (Ma et al., 2024). However, most existing methods are designed for post-event assessment, including expert inspections, ground sensors, and remote sensing (Li et al., 2021; Kucharczyk and Hugenholtz, 2021; Sarkar et al., 2023; Shi et al., 2022). While effective for characterizing observed damage, these approaches are reactive by nature and offer limited utility in pre-event planning, especially for sudden-onset events like earthquakes, where early awareness is crucial (Li et al., 2023). Traditional methods for pre-event simulation like scenario-based planning are useful (Ma et al., 2024;

Deierlein et al., 2020), but they need extensive domain expertise for region-specific models and often lack empirical validation by addressing human-centered factors.

Advances in large language models (LLMs) have shown promise in contextual simulation and complex reasoning across various domains (Li et al., 2024; Wang et al., 2024c). They are increasingly being viewed as potential world models—systems capable of learning to simulate and predict real-world scenarios (Wong et al., 2023; Hao et al., 2023). Through training on large-scale datasets that encode spatial, temporal, and causal relationships, LLMs have shown potential in learning representations of how the world works. For example, current research has demonstrated their ability to understand environment status (Hao et al., 2023), plan household activities (Xiang et al., 2023), and predict time-series events (Lee et al., 2025).

In disaster management, while LLMs are not yet widely applied as “world models,” researchers have explored their utility in tasks such as damage detection from satellite imagery (Zhang and Wang, 2024) or social media (Wang et al., 2024a), and emergency identification (Otal et al., 2024). However, key research gaps still remain. First, existing studies mainly use LLMs to analyze available textual or visual data for post-event assessment, not to simulate pre-event situations. Second, while LLMs have well-demonstrated reasoning abilities, effective disaster assessment requires integrating domain-specific knowledge and data fusion so that these models can accurately reason about potential disaster scenarios. To address these limitations, this study poses the fundamental question in the context of sudden-onset disasters: **Can LLMs simulate how humans perceive seismic risks before an event occurs?**

To answer this, we develop an LLM-based framework to simulate how humans perceive seismic risks, as illustrated in Figure 1. By integrating rich

pre-event contextual information, the LLMs are tasked with “reasoning” the likely severity of damage across spatial scales. Importantly, our study moves beyond theoretical simulations, grounding model evaluations in real-world events by testing on two actual earthquakes and comparing outputs against authoritative USGS reports. Our simulation demonstrates strong alignment with real reports at the zip code level, highlighting the potential of leveraging LLMs to improve pre-event planning.

2 Related Work

2.1 Seismic Hazard Simulation

Researchers have employed empirical, physics-based, and data-driven methods for simulating seismic disasters. Traditional approaches characterize earthquakes, such as the moment magnitude M_w (Moschetti et al., 2024), and then utilize empirical ground-motion prediction equations (GMPEs) (Moschetti et al., 2024; Iervolino, 2023) or physics-based simulations (Deierlein et al., 2020) to estimate site-specific shaking. While these simulations can capture complex local effects and rupture dynamics, they demand extensive data and computational resources. Recent data-driven methods have used machine learning to learn damage patterns from historical events, simulations, and remote sensing data (Cardellicchio et al., 2023; Yu et al., 2020)). These AI-driven techniques offer advantages in scalability and flexible feature integration but are dependent on high-quality labeled data, may struggle with generalization, and often present challenges in interpretability.

A major gap across these models, whether physics-based or data-driven, is the limited integration of human-perceived shaking into predictive frameworks. The USGS has developed the “Did You Feel It” (DYFI) system (Atkinson and Wald, 2007), a crowdsourced platform where individuals report the shaking intensity they experience after an earthquake. These reports are aggregated and converted into Modified Mercalli Intensity (MMI) (U.S. Geological Survey, 1989), which provide a human-centric, ground-truth reference for how seismic shaking is felt. While DYFI has been widely used to validate models or interpolate intensities after earthquakes, simulating human-centric perceived risks for pre-event could be important for developing predictive frameworks that anticipate not only physical shaking but also tangible impacts on communities.

2.2 LLM as World Models

Recent advance in LLMs such as GPT-4 (Achiam et al., 2023) and deepseek-R1 (Guo et al., 2025) has motivated researchers to leverage LLMs to solve complex tasks, such as reasoning (Li et al., 2025) and open-scene exploration (Wang et al., 2024b). In the context of disaster, researchers have leveraged LLMs to process multimodal information for vulnerability evaluation (Martelo et al., 2024), impact assessment (Taghian Dinani et al., 2023; Akinboyewa et al., 2024), information coordination (Yu and Wang, 2024), and recovery planning (White and Liptak, 2025).

Beyond conventional question-answering and analytical tasks, LLMs are being used as world models (Hao et al., 2023; Zhao et al., 2023) for complex simulation. While definitions of world models vary, their central concept involves leveraging forward reasoning to predict future states and dynamics in real-world scenarios. For example, they can produce large-scale simulations—an LLM-driven city simulation platform, OpenCity, runs tens of thousands of agents to emulate daily urban activities, successfully reproducing emergent patterns like mobility flows and residential segregation (Yan et al., 2024). In disaster scenarios, LLM-based world models have been explored in flooding forecasting (Wang et al., 2025) and evacuation planning (Hostetter et al., 2024).

Building on these advances, we propose leveraging LLMs as simulation tools to estimate how humans might perceive and report seismic risks before an event occurs. Our work addresses two critical gaps: first, the limited pre-event simulation of seismic risk; and second, the underexplored application of LLMs as world models for deriving human-centric insights.

3 Data and Methods

3.1 Framework Design

To simulate disaster impacts before events, we propose a novel framework that treats LLM as synthetic observers. This framework conceptualizes LLMs as “virtual sensors” capable of “perceiving” multimodal inputs and “reasoning” about disaster risks based on the MMI levels (see appendix A) that approximate human perception of shaking. As illustrated in Figure 1, we associate each sampled spatial location with a bundle of features, including seismic distance, site conditions, local building characteristics, and socioeconomic factors. In addi-

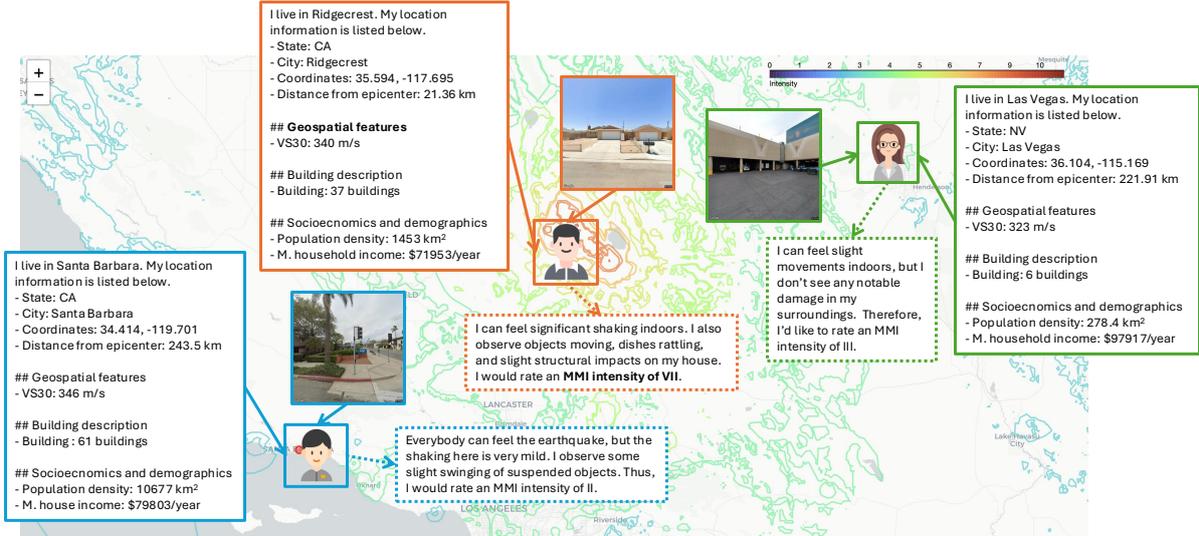


Figure 1: An illustration of LLM-simulated human-centric sensors.

tion, Google Street View imagery provides a first-person view of the built environment. Collectively, this feature set closely replicates the perceptual input available to disaster responders during an actual earthquake.

Formally, let each sample i be associated with a fused feature representation. We specifically select the following features that prior research has illustrated as critical predictors of seismic risk. (Frigerio et al., 2016; Kassem et al., 2020; Riedel et al., 2015; Mori et al., 2020).

$$\mathcal{X}_i = \{E_i, G_i, L_i, B_i, S_i, V_i\}$$

- E_i : Earthquake parameters (e.g., magnitude, epicenter distance, depth),
- G_i : Geospatial features (e.g., VS30),
- L_i : Location metadata (e.g., state, city, zip code),
- B_i : Building attributes (e.g., number, type, height, material),
- S_i : Socioeconomic indicators (e.g., population density, income),
- V_i : Street-level view (Google Street image).

The LLM acts as a reasoning function f_θ with parameters θ , generating both a reasoning trace and an MMI rating:

$$\hat{y}_i, e_i = f_\theta(\mathcal{X}_i), \quad \hat{y}_i \in \{\text{I, II, } \dots, \text{XII}\}$$

The full pipeline consists of five components, as shown in Figure 2: (1) spatial sampling, (2) data fusion, (3) prompt engineering, (4) experiment design, and (5) result analysis, which we specifically explain in the following sections.

3.2 Data Sampling

In Step 1 (Figure 2), we use a polygon-based GIS shapefile to define administrative zones (e.g., zip codes) and apply within-polygon stratified random sampling to ensure spatial representativeness across the study area. Let: $\mathcal{Z} = \{z_1, z_2, \dots, z_M\}$ represent the set of all zip code polygons. For each polygon z_j , we sample:

$$\mathcal{P}_j = \{p_{j1}, p_{j2}, \dots, p_{jn_j}\} \sim \text{Uniform}(z_j)$$

ensuring that points are uniformly drawn from within each polygon’s spatial boundary. We then sample 50 data points per zip code. This stratified random sampling strategy can help ensure spatial representativeness and mitigate biases associated with population density or urban–rural areas.

3.3 Data Fusion

In Step 2 (Figure 2), for each sampled point p_{ji} , we collect and assemble the feature set \mathcal{X}_{ji} from open-source and public datasets including USGS earthquake data, OpenStreetMap building data, American Community Survey (ACS) data, and Google Maps street views.

Earthquake data, site conditions, and location information. We source earthquake parameters, including moment magnitude, epicentral coordinates, and focal depth, from the USGS ShakeMap product (Wald et al., 2006) as E . To account for local site amplification and deamplification, we incorporate the USGS VS30 dataset (McPhillips et al., 2020) as the geospatial features G , a widely used proxy for near-surface geological effects.

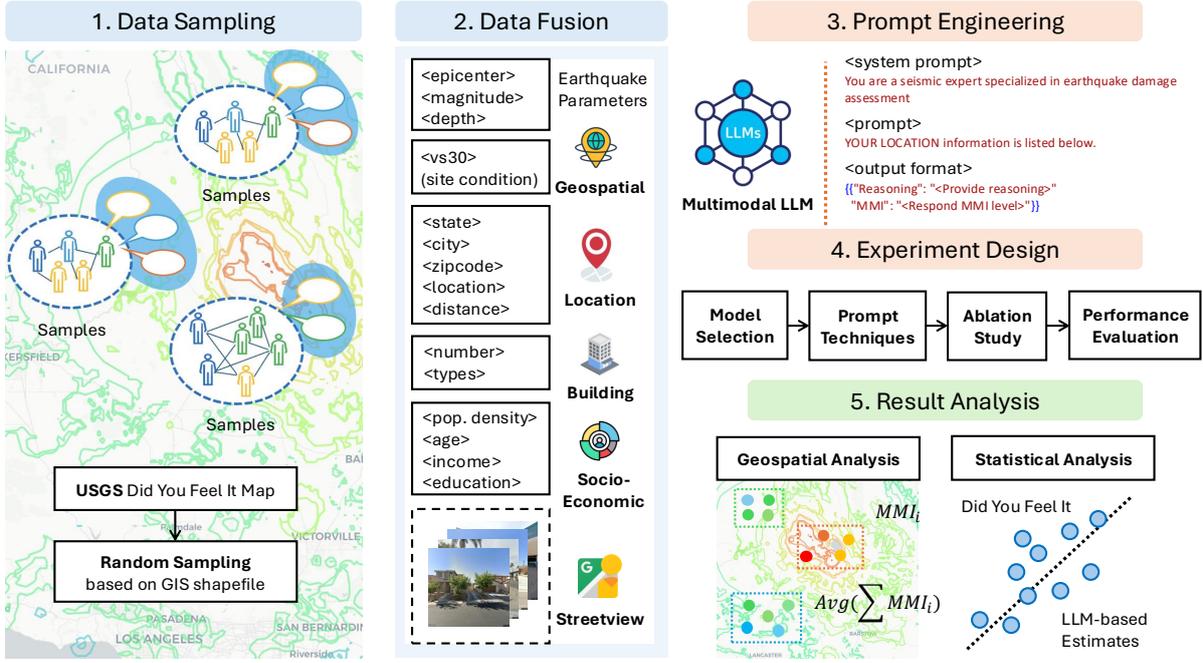


Figure 2: An illustration of the framework design.

We also incorporate location information, including latitude, longitude, state, county, zip code, and the distance from the epicenter for each sampled point. Additionally, we utilize MMI measurements from the USGS DYFI program as ground-truth labels, represented by y_j for each zip code.

Building description. We use OpenStreetMap (Ramm and Topf, 2010) (OSM) data to extract building features B , including the total number of buildings, type distribution, height range, and average height within a 100-meter radius of each sampled point. Meanwhile, we summarize the prevalence of major construction materials (e.g., concrete, masonry, timber, steel). These features can help characterize the distribution, physical attributes, and seismic-design status of buildings within the sensor’s surroundings.

Socioeconomic and demographic factors. We collect socioeconomic and demographic factors S at the Census Block Group (CBG) level from the American Community Survey (ACS) (USC, 2022). We spatially join the coordinates of each randomly selected sample to their corresponding CBG polygon and extract relevant ACS key indicators, including population, population density, urbanization ratio, 65- and over-age proportion, median household income, and higher education attainment rate. These variables imply a quantitative evaluation of population vulnerability in a disaster event.

Google Maps street imagery. We further incorporate Google Street View imagery V to enrich the environmental context at each sampled location. These street-level images capture fine-grained visual cues of the surrounding built environment—such as building, vegetation, curb conditions, and street density—that are difficult to numerically encode but essential for human-like visual reasoning. Using the Google Maps API (Google, 2025), we query the available image for each point in our sampling grid. This visual context allows the LLM to “see” the landscape as if it conducts a field visit.

3.4 Prompt Design

To guide the reasoning process of the LLM and ensure interpretability and consistency, we design a prompt template that mirrors the workflow of a seismic domain expert. The prompt follows a role-based instruction format in which the model is assigned the role of a seismic specialist responsible for evaluating earthquake damage using the MMI scale. The MMI scale provides a human-centric interpretation that consists of twelve levels describing the severity of earthquake shaking. The detailed descriptors of MMI are attached in Appendix A.

The prompt includes six distinct sections corresponding to the multimodal feature set \mathcal{X}_i introduced earlier: earthquake parameters E_i , geospatial features G_i , location metadata L_i , building descriptions B_i , socioeconomic context S_i , and street-

level view V_i . The model’s response is constrained to a JSON format with two fields: a free-text reasoning explanation and the predicted MMI category (see prompt template in Appendix B). We implement Chain-of-Thought (CoT) to instruct LLM output a detailed reasoning justifying its assessment before final MMI estimate.

3.5 Experiment Design

LLM Selection: We select both open- and closed-source LLMs to conduct the simulation. For open-source LLMs, we use models from Llama-3.2 and Qwen-2.5 family with various sizes, as listed in Table 1.

Prompting Techniques: We incorporate the following prompting techniques in our experiment:

- **In-Context Learning (ICL):** This helps LLM adapt to tasks by incorporating examples directly within the input prompt. ICL is applied by embedding a detailed MMI reference guide directly within the prompt.
- **Retrieval-Augmented Generation (RAG):** It combines information retrieval with text generation that allows LLMs to access external data sources at inference. We provide LLM with a set of multimodal features and the reported MMI within the prompt as the retrieved context to ground their MMI estimates.

Evaluation: The LLM generates a MMI prediction, $\hat{y}_{ji} = f_{\theta}(\mathcal{X}_{ji})$, for each sampled location i within a predefined administrative area j (e.g., zip code, county). These point-level predictions are then aggregated to compute an average predicted MMI for each area j . If area j contains n_j samples, its average predicted MMI $\hat{\bar{y}}_j$ is calculated as:

$$\hat{\bar{y}}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} \hat{y}_{ji}$$

Similarly, a corresponding ground-truth MMI value for area j , denoted as \bar{y}_j , is derived from USGS DYFI reports. To quantify the model’s predictive accuracy using these aggregated area-level values, we compute two metrics. First, the Root Mean Square Error (RMSE) is used:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{j=1}^N (\hat{\bar{y}}_j - \bar{y}_j)^2}$$

where N is the total number of administrative areas being evaluated (indexed by j). Second, we

calculate Pearson’s correlation coefficients r to assess the strength and direction of the association between the LLM-predicted MMI ($\hat{\bar{y}}_j$) and the ground-truth MMI (\bar{y}_j). These evaluations are performed at both zip code and county levels.

4 Experimental Results

We select two cases to demonstrate our proposed framework: (1) the 2014 Napa earthquake (magnitude 6.0) and (2) the 2019 Ridgecrest earthquake (magnitude 7.1), both of which occurred in California, U.S. (see details in Appendix C). For each case, we use the USGS DYFI reports as the ground-truth dataset (U.S. Geological Survey, 2014a, 2019a). Additionally, for each case, We compile feature sets for 50 sample points from each of the top 100 zip codes with the highest number of responses, resulting in 5,000 samples per event. Due to limitations in Google image availability for the Napa case, only 4,920 samples are retrieved. Comparisons with DYFI data are first conducted at the zip code level, where each of the 100 aggregated values represents the average of 50 simulated samples.

Figure 7 presents the spatial distribution of predicted MMI at the zip code level for the Napa and Ridgecrest earthquakes. This visualization highlights variations in simulated seismic impacts across geographical areas and among different LLMs. Based on the best-performing models (lowest RMSE: GPT-4.1-mini for the 2019 Ridgecrest earthquake and Qwen-2.5-32B for the 2014 Napa earthquake), we observe consistent geospatial patterns in both cases. Specifically, the simulations indicate elevated perceived risk near the epicenter (a red star mark in Figure 7), with diminishing simulated impact as distance increases. Moreover, the LLM-based predictions align well with DYFI reports from these two events: it is important to note that the Napa earthquake, despite its lower magnitude, led to more significant impacts.

The following sections are organized below. First, we evaluate the performance of the selected LLMs, comparing their accuracy using quantitative metrics, and examining the influence of model scaling and prompting strategies. Next, we conduct an input feature analysis to evaluate how different data modalities can impact predictive performance. Lastly, the output reasoning analysis explores the internal decision-making processes of the models, which identifies linguistic nuances that illustrate how LLMs interpret the inputs.

Model	Open Source	Napa				Ridgecrest			
		RMSE _Z ↓	Corr _Z ↑	RMSE _C ↓	Corr _C ↑	RMSE _Z ↓	Corr _Z ↑	RMSE _C ↓	Corr _C ↑
GPT-4o-2024-08-06	✗	2.43	0.77	2.37	0.88	1.97	0.75	1.91	0.77
GPT-4.1-mini	✗	2.56	0.61	2.48	0.67	0.92	0.64	0.77	0.76
Claude-3.5-haiku	✗	2.11	0.58	2.05	0.70	1.35	0.59	1.38	0.71
Llama-3.2-11B-VI	✓	3.19	0.44	3.05	0.86	3.22	0.33	3.22	0.27
Llama-3.2-90B-VI	✓	2.62	0.57	2.55	0.66	2.06	0.62	2.19	0.59
Qwen2.5-VL-3B	✓	3.63	0.29	3.59	0.15	3.88	0.01	4.08	-0.20
Qwen2.5-VL-7B	✓	1.79	0.43	1.68	0.70	1.53	0.05	1.59	-0.18
Qwen2.5-VL-32B	✓	1.59	0.70	1.56	0.79	0.99	0.71	0.96	0.80
Qwen2.5-VL-72B	✓	2.17	0.46	2.12	0.44	1.39	0.64	1.28	0.86

Table 1: Main experiment results, the best result for each metric is in bold.

4.1 Model Performance

Before the main experiment, we perform a data leakage test using two close-source models—Claude-3.5-haiku and GPT-4.1-mini—to show that our simulation is free of data leakage issues (see Appendix D). The main experiment results are shown in Table 1, from which we draw the following findings:

open-source models. Table 1 shows a clear trend that closed-source LLMs consistently outperform open-source counterparts, achieving the best results in 6 out of 8 cases. This suggests that commercial closed-source models possess stronger geospatial reasoning capabilities and align more closely with human judgment in disaster sensing tasks. Among the open-source models, Qwen-2.5-32B obtains the top results with lowest RMSE for the Napa case.

RMSE and correlation can be inconsistent. Another noteworthy observation from Table 1 is that the two evaluation metrics—correlation and RMSE—do not always align. For instance, Llama-3.2-11B shows a high correlation but a poor RMSE at the county level for the Napa case. This discrepancy arises because correlation captures the model’s ability to predict relative ordering of seismic impacts, whereas RMSE reflects the absolute prediction errors. Thus, even when models effectively estimate the relevant severity (lower RMSE), they may not correctly distinguish between higher- or lower-impact areas (lower correlation).

We then evaluate the influence of incorporating demonstration methods with RAG and ICL on the simulation, as shown in Figure 3. These demonstration techniques can enhance the capacity of LLMs to generalize from provided information. In particular, the enhanced models consistently show higher performance, illustrated by decreasing $RMSE_Z$. It is also well-noted that even limited demonstrations can greatly improve model predictions.

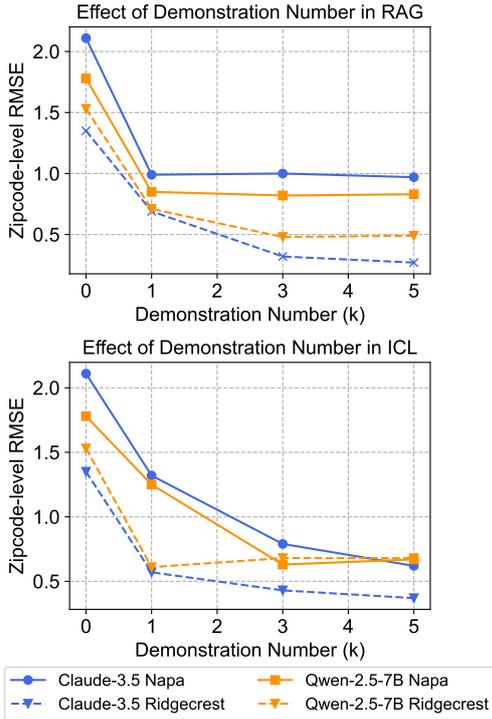


Figure 3: Demonstration analysis on RAG and ICL.

LLMs deliver promising simulation results. As shown in Table 1, all of the best correlation scores across datasets and area levels exceed 0.7, indicating a strong alignment between the predicted outcomes and the ground truth impact labels. These high correlation values suggest that LLMs hold significant potential for effectively simulating human-perceived risks in disaster scenarios.

Closed-source LLMs generally outperform

4.2 Input Feature Analysis

In this section, we examine how input features beyond earthquake and location information affect LLMs’ simulation performance. We conduct experiments using Claude-3.5-haiku and Qwen-2.5-7B, with the results presented in Figure 4. Interestingly, we find that only street view information contributes to improved simulation perfor-

mance. In contrast, removing any of the other three features alone—geospatial, building, or socioeconomic data—leads to an increase in zip code-level RMSE.

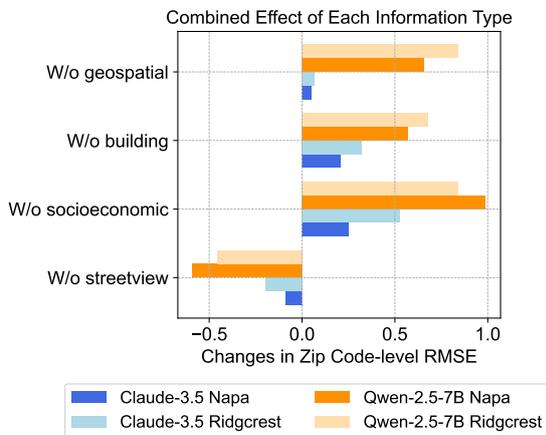


Figure 4: Input information feature analysis results.

We attribute this performance degradation to several potential factors: (1) limitations of LLMs in processing numerical information as they could complicate LLM’s reasoning process (Yin et al., 2024; Bodensohn et al., 2025); (2) the lack of domain-specific knowledge required to interpret geospatial, building, and community-related data (Gao et al., 2024); (3) the inherent limitations of the self-attention mechanism in capturing spatial adjacency and performing geometric reasoning (Requeima et al., 2024).

4.3 Output Reasoning Analysis

To complement our quantitative evaluations, we examine how GPT-4.1-mini and Qwen-2.5-32B (lowest RMSE) reason when predicting MMI values. We summarize our findings as below:

LLMs capture seismic attenuation but underutilize local site conditions. As shown in Figure 5, both models display a clear negative correlation between epicentral distance and predicted MMI, most notably in Qwen2.5-32B’s Napa earthquake predictions, which indicates that LLMs have internalized the concept of seismic attenuation. However, the relationship between VS30 values (a proxy for local ground conditions) and MMI is weak across both models. High MMI values occur almost exclusively near the epicenter, suggesting limited sensitivity to local site effects.

LLMs use distinct lexical cues for MMI reasoning across multimodal inputs. Figure 6 presents a

taxonomy of language used by both models across three perspectives: buildings, socioeconomic context, and street-level imagery. Unigram analysis reveals that GPT-4.1-mini and Qwen2.5-32B employ different vocabularies and reasoning styles. For the building assessment, GPT-4.1-mini adopts a descriptive and hedged reasoning style. At low MMI levels, it uses terms like “enhance,” “robust,” and “improve,” while shifting to “suffer” and “detached” at higher levels. Qwen2.5-32B relies on more technical terms, such as “compliance” and “stringent” at low levels, and “crack” and “susceptibility” at higher ones.

Socioeconomic reasoning diverges in focus and tone. GPT-4.1-mini links lower MMI levels to terms like “urbanized” and “welleducated,” and higher levels to “industrial” and “heavy,” occasionally incorporating cautious language such as “possibly” or “suggest.” Qwen2.5-32B emphasizes systemic vulnerability, shifting from “limited” and “stable” to “cascade” and “amplification” as predicted MMI increases.

Visual reasoning contrasts environmental vs. structural emphasis. GPT-4.1-mini references broader environmental cues, from “forest” and “agricultural” at lower levels to “catastrophic” at higher ones. In contrast, Qwen2.5-32B focuses on structural compliance, mentioning “firm,” “code,” and “reinforced” in a progression of seismic risk.

5 Discussion

Based on our experimental result analysis, we conclude the following insights and opportunities for future study:

- LLMs can effectively simulate human-centric seismic risks, showing strong alignment with real-world USGS DYFI reports. This can help develop pre-event impact assessment. One promising direction is to explore broader application in LLM-based simulations for disaster responses.
- Simulation performance depends on model type, size, and input information: closed-source and larger models generally perform better; techniques like RAG and ICL help; street-level imagery boosts accuracy, while structured data may hinder it due to modality alignment limitations. Future works can explore more effective prompting strategies and reasoning structures to further improve the simulation performance.

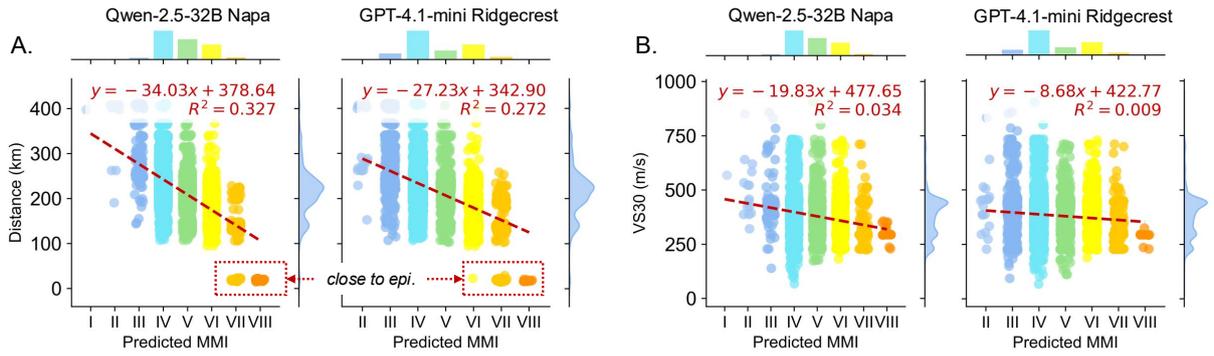


Figure 5: Output reasoning analysis in terms of (A) distance (where the x-axis is the predicted MMI, and y-axis is the distance from the epicenter (km)) and (B) VS30 (where the x-axis is the predicted MMI, and the y-axis the local site condition represented by VS30 (m/s)).

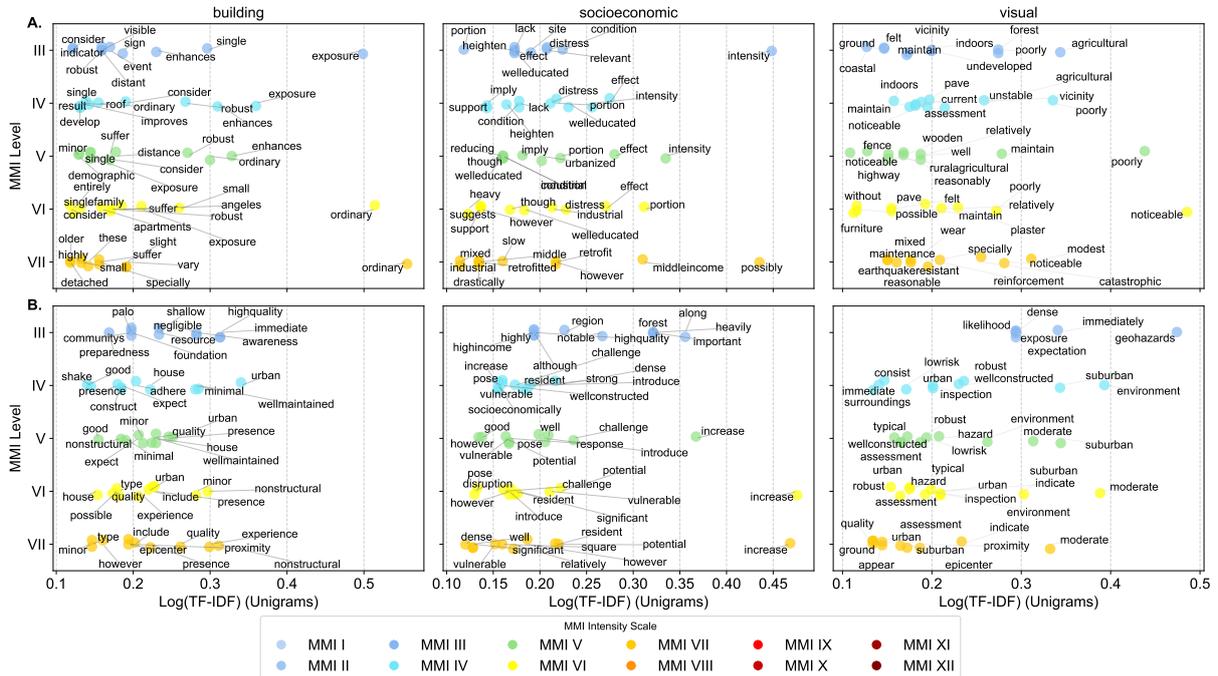


Figure 6: Output reasoning analysis in terms of different input features with (A) Qwen-2.5-32B for the Napa case, and (B) GPT-4.1-mini for the Ridgecrest case. The x-axis is the log(TF-IDF), while the y-axis is the predicted MMI.

- LLMs show diverse reasoning styles and strong practical value, as model-specific rhetorical patterns reflect architectural and data differences. These differences suggest the potential impact of training data and model architecture on LLMs' reasoning styles. It would be meaningful to further explore the reasoning mechanisms of LLMs when used as world models.

Our study has significant practical implications. Leveraging LLMs and open-source data to simulate seismic risk enables authorities to proactively assess potential disaster impacts. Integrating our framework into early-warning systems can also help identify vulnerable communities and enhance

disaster response.

6 Conclusions

Our study demonstrates the potential of using LLMs to simulate seismic risk before an earthquake occurs. The alignment between model predictions and real-world reports highlights the importance of multimodal inputs and advanced LLM techniques like RAG and ICL in simulation settings. Moreover, our findings reveal that simulation performance can vary across LLMs and input features. These results make a meaningful contribution to advancing data-driven, human-centric simulation with LLMs for real-world scenarios.

7 Limitations

Several limitations warrant further consideration. First, our experiments focus only on two earthquake cases—the 2014 Napa and 2019 Ridgecrest events—which may not fully represent global variations in seismic hazards, urban densities, and construction practices. Therefore, generalizability requires additional validation.

Second, despite employing stratified random sampling within zip codes, potential biases remain due to gaps in data availability, such as incomplete Google Street View coverage or inconsistencies in socioeconomic and building datasets. This may lead to underrepresentation of certain neighborhoods or misalignment with high-risk zones.

Third, although our framework integrates diverse heterogeneous data, we do not conduct detailed feature selection or examine individual parameters (e.g., housing age, infrastructure proximity). Consequently, interpreting model behavior at a granular level requires further investigation.

8 Ethics

Responsible data use. Our research relies exclusively on publicly available and anonymized datasets, including geospatial, demographic, and imagery data from USGS, OpenStreetMap, and Google Street View. All analyses are conducted at aggregated administrative levels (zip code and county), minimizing risks to individual privacy.

Deployment guidance. It is important to recognize that LLM-based simulations cannot fully capture the complexity and diversity of individual experiences in real disaster scenarios. Therefore, model predictions should be viewed as supportive tools rather than replacements for expert judgment, and used responsibly with validated empirical assessments in practical decision-making contexts.

References

2022. [American community survey \(acs\)](#).

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. [Gpt-4 technical report](#). *ArXiv preprint*, abs/2303.08774.

Temitope Akinboye, Huan Ning, M Naser Lessani, and Zhenlong Li. 2024. Automated floodwater depth estimation using large multimodal model for

rapid flood mapping. *Computational Urban Science*, 4(1):12. 615 616

Gail M. Atkinson and David J. Wald. 2007. “did you feel it?” intensity data: A surprisingly good measure of earthquake ground motion. *Seismological Research Letters*, 78(3):362–368. 617 618 619 620

Jan-Micha Bodensohn, Ulf Brackmann, Liane Vogel, Anupam Sanghi, and Carsten Binnig. 2025. Unveiling challenges for llms in enterprise data engineering. *arXiv preprint arXiv:2504.10950*. 621 622 623 624

Angelo Cardellicchio, Sergio Ruggieri, Valeria Leggieri, and Giuseppina Uva. 2023. A machine learning framework to estimate a simple seismic vulnerability index from a photograph: the vulma project. *Procedia Structural Integrity*, 44:1956–1963. XIX ANIDIS Conference, Seismic Engineering in Italy. 625 626 627 628 629 630

Gregory G Deierlein, Frank McKenna, Adam Zsarnóczay, Tracy Kijewski-Correa, Ahsan Kareem, Wael Elhaddad, Laura Lowes, Matthew J Schoetler, and Sanjay Govindjee. 2020. A cloud-enabled application framework for simulating regional-scale impacts of natural hazards on the built environment. *Frontiers in Built Environment*, 6:558706. 631 632 633 634 635 636 637

Ivan Frigerio, Stefania Ventura, Daniele Strigaro, Matteo Mattavelli, Mattia De Amicis, Silvia Mugnano, and Mario Boffi. 2016. A gis-based approach to identify the spatial variability of social vulnerability to seismic hazard in italy. *Applied geography*, 74:12–22. 638 639 640 641 642 643

Yanjun Gao, Skatje Myers, Shan Chen, Dmitriy Dligach, Timothy A Miller, Danielle Bitterman, Matthew Churpek, and Majid Afshar. 2024. When raw data prevails: Are large language model embeddings effective in numerical data representation for medical machine learning applications? *arXiv preprint arXiv:2408.11854*. 644 645 646 647 648 649 650

Google. 2025. Google Maps Platform. <https://developers.google.com/maps>. Accessed: 2025-05-19. 651 652 653

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*. 654 655 656 657 658 659

Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*. 660 661 662 663

Haley Hostetter, MZ Naser, Xinyan Huang, and John Gales. 2024. Large language models in fire engineering: An examination of technical questions against domain knowledge. *arXiv preprint arXiv:2403.04795*. 664 665 666 667 668

669	Junio Iervolino. 2023. Implications of gmpe’s structure for multi-site seismic hazard. <i>Soil Dynamics and Earthquake Engineering</i> , 172:108022.	725
670		726
671		727
672	Rebecca Louise Jones, Debarati Guha-Sapir, and Sandy Tubeuf. 2022. Human and economic impacts of natural disasters: can we trust the global data? <i>Scientific data</i> , 9(1):572.	728
673		729
674		730
675		731
676	Moustafa Moufid Kassem, Fadzli Mohamed Nazri, and Ehsan Noroozinejad Farsangi. 2020. The seismic vulnerability assessment methodologies: A state-of-the-art review. <i>Ain Shams Engineering Journal</i> , 11(4):849–864.	732
677		733
678		734
679		735
680		736
681	Maja Kucharczyk and Chris H Hugenholtz. 2021. Remote sensing of natural hazard-related disasters with small drones: Global trends, biases, and research opportunities. <i>Remote Sensing of Environment</i> , 264:112577.	737
682		738
683		739
684		740
685		741
686	Geon Lee, Wenchao Yu, Kijung Shin, Wei Cheng, and Haifeng Chen. 2025. Timecap: Learning to contextualize, augment, and predict time series events with large language model agents. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 18082–18090.	742
687		743
688		744
689		745
690		746
691		747
692	Bowen Li, Zhaoyu Li, Qiwei Du, Jinqi Luo, Wenshan Wang, Yaqi Xie, Simon Stepputtis, Chen Wang, Kattia Sycara, Pradeep Ravikumar, and 1 others. 2024. Logicity: Advancing neuro-symbolic ai with abstract urban simulation. <i>Advances in Neural Information Processing Systems</i> , 37:69840–69864.	748
693		749
694		750
695		751
696		752
697		753
698	Lingyao Li, Michelle Bensi, and Gregory Baecher. 2023. Exploring the potential of social media crowdsourcing for post-earthquake damage assessment. <i>International Journal of Disaster Risk Reduction</i> , 98:104062.	754
699		755
700		756
701		757
702		758
703	Lingyao Li, Michelle Bensi, Qingbin Cui, Gregory B Baecher, and You Huang. 2021. Social media crowdsourcing for rapid damage assessment following a sudden-onset natural hazard event. <i>International Journal of Information Management</i> , 60:102378.	759
704		760
705		761
706		762
707		763
708	Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, and 1 others. 2025. From system 1 to system 2: A survey of reasoning large language models. <i>arXiv preprint arXiv:2502.17419</i> .	764
709		765
710		766
711		767
712		768
713		769
714	Zihui Ma, Lingyao Li, Yujie Mao, Yu Wang, Olivia Grace Patsy, Michelle T Bensi, Libby Hemphill, and Gregory B Baecher. 2024. Surveying the use of social media data and natural language processing techniques to investigate natural disasters. <i>Natural Hazards Review</i> , 25(4):03124003.	770
715		771
716		772
717		773
718		774
719		775
720	Rafaela Martelo, Kimia Ahmadiyehyazdi, and Ruoqian Wang. 2024. Towards democratized flood risk management: An advanced ai assistant enabled by gpt-4 for enhanced interpretability and public engagement. <i>arXiv preprint arXiv:2403.03188</i> .	776
721		777
722		778
723		779
724		780
	Devin F McPhillips, Julie A Herrick, Sean Ahdi, Alan K Yong, and Scott Haefner. 2020. Updated compilation of vs30 data for the united states. (<i>No Title</i>).	
	Federico Mori, Amerigo Mendicelli, Massimiliano Moscatelli, Gino Romagnoli, Edoardo Peronace, and Giuseppe Naso. 2020. A new vs30 map for italy based on the seismic microzonation dataset. <i>Engineering Geology</i> , 275:105745.	
	Morgan P Moschetti, Brad T Aagaard, Sean K Ahdi, Jason Altekruze, Oliver S Boyd, Arthur D Frankel, Julie Herrick, Mark D Petersen, Peter M Powers, Sanaz Rezaeian, and 1 others. 2024. The 2023 us national seismic hazard model: Ground-motion characterization for the conterminous united states. <i>Earthquake Spectra</i> , 40(2):1158–1190.	
	Hakan T Otal, Eric Stern, and M Abdullah Canbaz. 2024. Llm-assisted crisis management: Building advanced llm platforms for effective emergency response and public collaboration. In <i>2024 IEEE Conference on Artificial Intelligence (CAI)</i> , pages 851–859. IEEE.	
	Frederik Ramm and Jochen Topf. 2010. <i>OpenStreetMap: Die freie Weltkarte nutzen und mitgestalten</i> . Lehmanns Media.	
	James Requeima, John Bronskill, Dami Choi, Richard Turner, and David K Duvenaud. 2024. Llm processes: Numerical predictive distributions conditioned on natural language. <i>Advances in Neural Information Processing Systems</i> , 37:109609–109671.	
	Ismaël Riedel, Philippe Guéguen, Mauro Dalla Mura, Erwan Pathier, Thomas Leduc, and Jocelyn Chanussot. 2015. Seismic vulnerability assessment of urban environments in moderate-to-low seismic hazard regions using association rule learning and support vector machine methods. <i>Natural hazards</i> , 76:1111–1141.	
	Argho Sarkar, Tashnim Chowdhury, Robin Roberson Murphy, Aryya Gangopadhyay, and Maryam Rah-nemoonfar. 2023. Sam-vqa: Supervised attention-based visual question answering model for post-disaster damage assessment on remote sensing imagery. <i>IEEE Transactions on Geoscience and Remote Sensing</i> , 61:1–16.	
	Kaize Shi, Xueping Peng, Hao Lu, Yifan Zhu, and Zhen-dong Niu. 2022. Application of social sensors in natural disasters emergency management: a review. <i>IEEE Transactions on Computational Social Systems</i> , 10(6):3143–3158.	
	Soudabeh Taghian Dinani, Doina Caragea, and Nikesh Gyawali. 2023. Disaster tweet classification using fine-tuned deep learning models versus zero and few-shot large language models. In <i>International Conference on Data Management Technologies and Applications</i> , pages 73–94. Springer.	
	U.S. Geological Survey. 1989. The modified mercalli intensity (mmi) scale. https://pubs.usgs.gov .	

781	gov/gip/earthq4/severitygip.html . Accessed:	and Joshua B Tenenbaum. 2023. From word mod-	835
782	2024-05-18.	els to world models: Translating from natural lan-	836
783	U.S. Geological Survey. 2014a. Did you	guage to the probabilistic language of thought. <i>arXiv</i>	837
784	feel it? – community internet intensity	<i>preprint arXiv:2306.12672</i> .	838
785	map for 2014 napa earthquake. https:	839	
786	://earthquake.usgs.gov/earthquakes/	840	
787	eventpage/nc72282711/dyfi/responses . Ac-	841	
788	cessed: 2024-05-18.	842	
789	U.S. Geological Survey. 2014b. M 6.0 -	843	
790	south napa earthquake, california. https:	844	
791	://earthquake.usgs.gov/earthquakes/	845	
792	eventpage/nc72282711/executive . Accessed:	846	
793	2024-05-18.	847	
794	U.S. Geological Survey. 2019a. Did you	848	
795	feel it? – community internet inten-	849	
796	sity map for 2019 ridgecrest earthquake.	850	
797	https://earthquake.usgs.gov/earthquakes/	851	
798	eventpage/ci38457511/dyfi/intensity . Ac-	852	
799	cessed: 2024-05-18.	853	
800	U.S. Geological Survey. 2019b. M 7.1	854	
801	- 2019 ridgecrest earthquake sequence.	855	
802	https://earthquake.usgs.gov/earthquakes/	856	
803	eventpage/ci38457511/executive . Accessed:		
804	2024-05-18.	857	
805	David J Wald, Bruce C Worden, Vincent Quitoriano, and	858	
806	Kris L Pankow. 2006. Shakemap® manual. <i>Techni-</i>	859	
807	<i>cal Manual, users guide, and software guide Version</i> .	860	
808	Chenguang Wang, Davis Engler, Xuechun Li, James	861	
809	Hou, David J Wald, Kishor Jaiswal, and Susu Xu.	862	
810	2024a. Near-real-time earthquake-induced fatality es-		
811	timation using crowdsourced data and large-language	863	
812	models. <i>International Journal of Disaster Risk Re-</i>	864	
813	<i>duction</i> , 111:104680.	865	
814	Gelan Wang, Yu Liu, Shukai Liu, Ling Zhang, and	866	
815	Liqun Yang. 2025. Remflow: Rag-enhanced multi-	867	
816	factor rainfall flooding warning in sponge airports	868	
817	via large language model. <i>International Journal of</i>	869	
818	<i>Machine Learning and Cybernetics</i> , pages 1–21.	870	
819	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	871	
820	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,		
821	Xu Chen, Yankai Lin, and 1 others. 2024b. A survey		
822	on large language model based autonomous agents.		
823	<i>Frontiers of Computer Science</i> , 18(6):186345.		
824	Yue Wang, Tianfan Fu, Yinlong Xu, Zihan Ma, Hongxia		
825	Xu, Bang Du, Yingzhou Lu, Honghao Gao, Jian Wu,		
826	and Jintai Chen. 2024c. Twin-gpt: digital twins for		
827	clinical trials via large language model. <i>ACM Trans-</i>		
828	<i>actions on Multimedia Computing, Communications</i>		
829	<i>and Applications</i> .		
830	Gwen White and Sadie Liptak. 2025. Small business		
831	continuity and disaster recovery plans using ai and		
832	chatgpt.		
833	Lionel Wong, Gabriel Grand, Alexander K Lew, Noah D		
834	Goodman, Vikash K Mansinghka, Jacob Andreas,		

A Modified Mercalli Intensity (MMI) Scale

Table 2 shows the MMI scale used to support the classification of seismic risks in our study.

MMI Level	Description
I	Not felt except by a very few under especially favorable conditions.
II	Felt only by a few persons at rest, especially on upper floors of buildings. Delicately suspended objects may swing.
III	Felt quite noticeably by persons indoors, especially on upper floors of buildings. Many people do not recognize it as an earthquake. Standing motor cars may rock slightly. Vibration similar to the passing of a truck. Duration estimated.
IV	Felt indoors by many, outdoors by few during the day. At night, some awakened. Dishes, windows, doors disturbed; walls make cracking sound. Sensation like heavy truck striking building. Standing motor cars rocked noticeably.
V	Felt by nearly everyone; many awakened. Some dishes, windows broken. Unstable objects overturned. Pendulum clocks may stop.
VI	Felt by all, many frightened. Some heavy furniture moved; a few instances of fallen plaster. Damage slight.
VII	Damage negligible in buildings of good design and construction; slight to moderate in well-built ordinary structures; considerable damage in poorly built or badly designed structures; some chimneys broken.
VIII	Damage slight in specially designed structures; considerable damage in ordinary substantial buildings with partial collapse. Damage great in poorly built structures. Fall of chimneys, factory stacks, columns, monuments, walls. Heavy furniture overturned.
IX	Damage considerable in specially designed structures; well-designed frame structures thrown out of plumb. Damage great in substantial buildings, with partial collapse. Buildings shifted off foundations.
X	Some well-built wooden structures destroyed; most masonry and frame structures destroyed with foundations. Rail bent.
XI	Few, if any (masonry) structures remain standing. Bridges destroyed. Rails bent greatly.
XII	Damage total. Lines of sight and level are distorted. Objects thrown into the air.

Table 2: Description of the Modified Mercalli Intensity (MMI) scale (U.S. Geological Survey, 1989).

B Prompt Template Design

```

1 SYSTEM_PROMPT = ""
2 You are a seismic expert specialized in earthquake damage assessment and
  disaster response. You analyze earthquake data, local conditions, and
  building characteristics to provide damage assessments using the Modified
  Mercalli Intensity (MMI) scale.
3 ""

```

```

1 EARTHQUAKE_PROMPT = ""
2 The earthquake happened date is 2025-06-01.
3
4 Here is the EARTHQUAKE information.
5 - Epicenter: {eq_place}
6 - Coordinates: {eq_lat}, {eq_lng}
7 - Magnitude: {eq_magnitude} mw
8 - Depth: {eq_depth} km
9
10 YOUR LOCATION information is listed below.
11 - State: {state}
12 - City: {city}
13 - Zipcode: {zipcode}
14 - Coordinates: {lat}, {lng}
15 - Distance from epicenter: {distance} km
16
17 ## Geospatial features in YOUR LOCATION
18 - VS30 at your location: {vs30} m/s
19 (VS30 represents the time-averaged shear-wave velocity (VS) to a depth of 30
  meters, which is a key index to account for seismic site conditions)
20
21 ## Building Description in YOUR LOCATION (within a 100-meter radius)
22 - Building description: {building}
23

```

```

24 ## Community Socioeconomics and Demographics in YOUR LOCATION (at Census Block
    Group level)
25 - Population density: {population_density} people per square km
26 - Urban population percentage: {urban_population_pct}%
27 - Over 65 percentage: {over_65_rate}%
28 - Median household income: ${median_household_income}/year
29 - Education (bachelor's or higher): {education}%
30
31 ## Visual Context in YOUR LOCATION
32 The image provided shows your surrounding environment and infrastructure.
33
34 Based on the information provided, ASSESS the potential earthquake damage level
    using the Modified Mercalli Intensity (MMI) scale.
35 1. Identify the damage level.
36 2. Explain your reasoning by addressing the following factors and considering
    the visual context.
37 - Distance to the epicenter and earthquake magnitude
38 - Geospatial features
39 - Infrastructure quality and building characteristics
40 - Population density and socioeconomic vulnerabilities
41 - Visual image of surroundings
42
43 The following is an abbreviated description of the 12 levels of Modified
    Mercalli intensity. {MMI Scale}
44
45 Output the result in JSON format:
46 {{
47     "Reasoning": "<Provide reasoning>"
48     "MMI": "<Respond MMI level>",
49 }}
50 ""

```

C Earthquake Scenarios

2014 Napa Earthquake (U.S. Geological Survey, 2014b). On August 24, 2014, a magnitude 6.0 earthquake struck near Napa, California, causing significant structural damage despite its moderate magnitude. Approximately 613 buildings were tagged for various degrees of structural integrity concerns, including fractures, road cracks, and damage to wine storage facilities. The earthquake resulted in one death and nearly 200 injuries.

2019 Ridgecrest Earthquake (U.S. Geological Survey, 2019b). The Ridgecrest earthquake occurred on July 6, 2019, with a magnitude of 7.1, significantly larger than the Napa event but with fewer human casualties. The quake damaged around 50 homes, caused gas leaks and road cracks, and triggered fires in residential properties. Significant infrastructural damage occurred at the Naval Air Station, and widespread power outages were reported.

D Data Leakage Test

To further assess the potential for data leakage in the LLMs used in our experiments, we conduct a leakage test on Claude-3.5-haiku and GPT-4.1-mini, the two best-performing models in our main study. Specifically, we remove city and state names from the prompt—two elements most likely to serve as shortcuts for the models to associate with MMI levels and potentially memorize. As shown in Table 3, the removal of location information does not significantly affect the models’ simulation performance. It is fair to rule out the possibility of data leakage in our main results.

Model	Earthquake prompt				Earthquake prompt w/o location			
	RMSE_Z	Corr_Z	RMSE_C	Corr_C	RMSE_Z	Corr_Z	RMSE_C	Corr_C
claude-3-5-haiku	2.11	0.58	2.05	0.70	2.35	0.38	2.26	0.62
gpt-4.1-mini	2.56	0.61	2.48	0.67	2.67	0.62	2.58	0.73

Table 3: Experiment results on data leakage test.

Model Family	Model Size	Napa RMSE	Ridgecrest RMSE
Llama-3.2	11B	3.19	3.22
	90B	2.62	2.06
Qwen-2.5	3B	3.63	3.88
	7B	1.79	1.53
	72B	2.17	1.39

Table 4: Scaling law analysis: Zipcode-level RMSE across model sizes (in billions of parameters).

E Spatial Distribution of Predicted MMI by LLMs

Figure 7 shows the spatial distribution of predicted MMI for the 2014 Napa earthquake and the 2019 Ridgecrest earthquake at zip code level from different LLMs.

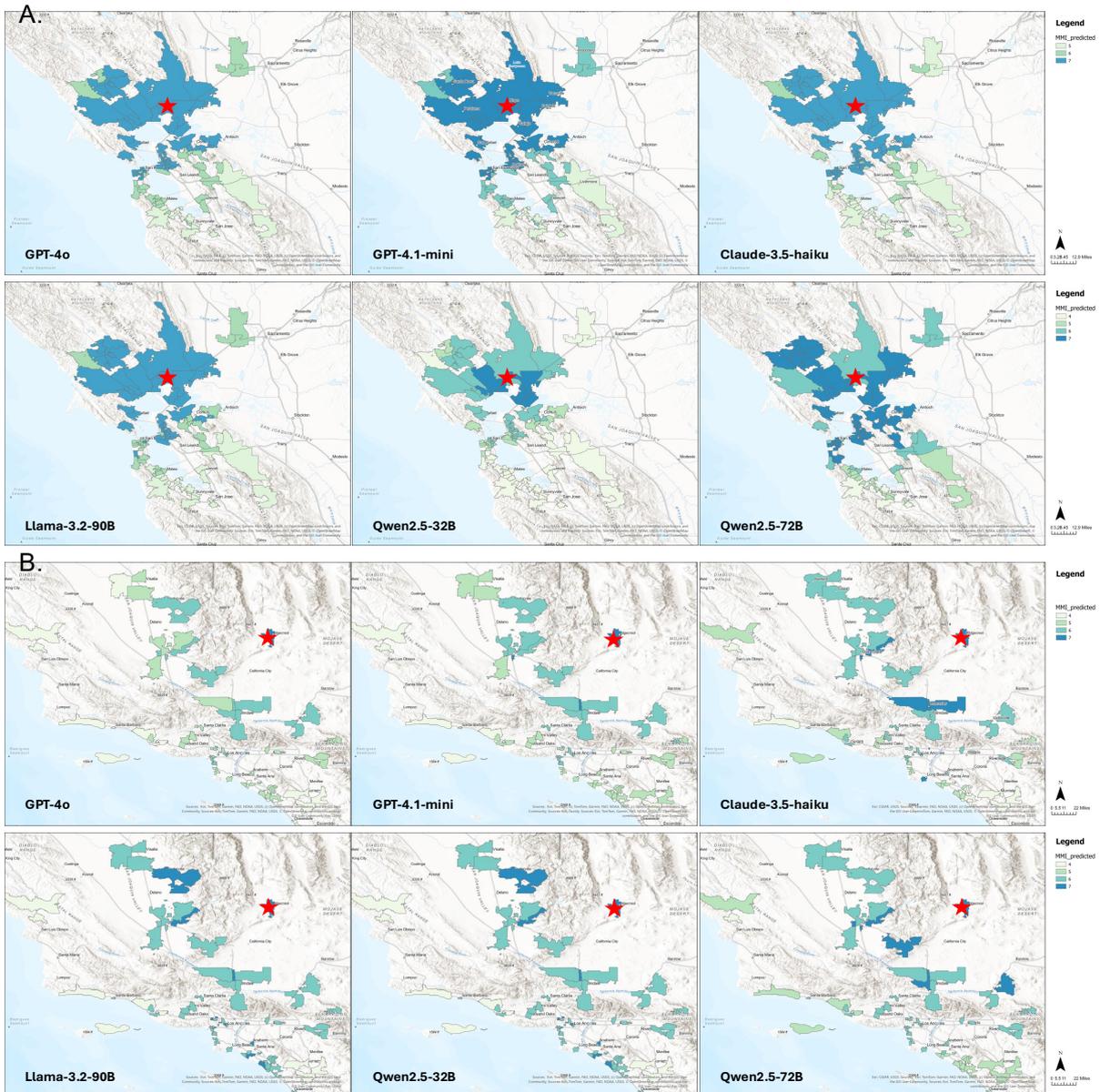


Figure 7: Spatial distribution of predicted MMI at zip code level: (A) The 2014 Napa earthquake, and (B) the 2019 Ridgecrest earthquake based on GPT4o, GPT-4.1-mini, Claude-3.5-haiku, Llama-3.2-90B, Qwen-2.5-32B, and Qwen-2.5-72B. The red star sign shows the epicenter of the earthquake. These maps compare outputs from different LLMs, showing consistent geospatial patterns with high-intensity predictions concentrated near epicenters

F Scaling Law Analysis

961

To investigate the impact of LLM scaling laws on simulation performance, we conduct a scaling analysis using the Llama-3.2 and Qwen-2.5 model families. As shown in Table 4, we observe that performance in simulation generally improves with model size, excepting Qwen-2.5-72B. This indicates that the simulation performance tend to be strengthened with a larger size of LLMs.

962

963

964

965

G Samples of LLM Output

966

Samples of LLM Reasoning Outputs with Different MMI Predictions

GPT-4.1-mini | 2019 Ridgecrest | Predicted MMI: III

The earthquake has a high magnitude of 7.1 Mw, which is capable of causing significant damage near the epicenter. However, the location in question, San Diego, CA, is approximately 339 km away from the epicenter in Ridgecrest, CA. At this distance, seismic waves would have attenuated considerably, reducing the intensity of shaking experienced. The VS30 value of 417 m/s indicates a site with relatively soft soil conditions, which can amplify shaking somewhat, but given the large distance, this effect is limited. The building stock consists mostly of general buildings and schools, likely built to modern California seismic codes, which enhances resilience. The population density is moderate, with a well-educated and relatively affluent community, reducing socioeconomic vulnerabilities. The visual context shows a well-maintained urban street with no visible signs of damage or distress, supporting the assessment of low shaking impact. Considering all these factors, the expected Modified Mercalli Intensity (MMI) level in San Diego would be around III, where shaking is noticeable indoors but causes no damage.

GPT-4.1-mini | 2019 Ridgecrest | Predicted MMI: IV

The earthquake has a high magnitude of 7.1 Mw, which is capable of causing significant damage near the epicenter. However, the location in Ventura, CA is approximately 221 km away from the epicenter in Ridgecrest, CA, which greatly reduces the intensity of shaking experienced. The VS30 value of 391 m/s indicates a site with relatively soft soil conditions, which can amplify shaking somewhat, but the distance still plays a dominant role in reducing impact. The building information is not available, but the visual context shows modern infrastructure such as a highway bridge and vehicles, suggesting reasonably good construction standards typical of California. The population density is low (65 people/km²) with a moderate urban population percentage (25.43 %), and socioeconomic indicators such as median income and education level are moderate, implying some resilience but also potential vulnerabilities in older or poorly maintained structures. Given the distance, moderate site conditions, and infrastructure, the expected shaking intensity would be moderate with little to no structural damage. This corresponds to an MMI level around IV to V, where shaking is felt by many but damage is minimal or limited to minor non-structural effects.

GPT-4.1-mini | 2019 Ridgecrest | Predicted MMI: V

The earthquake had a high magnitude of 7.1 Mw, which is capable of causing severe damage near the epicenter. However, the location in question, Apple Valley, CA, is approximately 152.25 km away from the epicenter, which significantly reduces the intensity of shaking experienced. The VS30 value of 607 m/s indicates relatively firm ground conditions, which typically reduce amplification of seismic waves. The visual context shows a rural, sparsely populated area with no visible dense infrastructure or tall buildings, suggesting limited exposure to structural damage. The population density is low (18.73 people per square km), and the urban population is only about 32.65 %, further reducing potential damage and casualties. Socioeconomic factors such as median income and education level are moderate, but given the lack of detailed building information and the rural setting, the risk of significant structural damage is low. Considering these factors, the shaking would likely be felt by many but cause only slight damage, consistent with MMI level VI.

Qwen-2.5-32B | 2014 Napa | Predicted MMI: VI

The location is 19.27 km from the epicenter of a 6.0 Mw earthquake, indicating strong shaking. The VS30 value of 344 m/s suggests stiff soil, reducing amplification but not eliminating damage due to proximity. The presence of general buildings and houses, along with a rural setting and a high percentage of older residents, points to moderate to strong damage. The visual context shows a well-maintained residential area, which may indicate better-quality buildings, but this is not definitive.

Qwen-2.5-32B | 2014 Napa | Predicted MMI: VII

The earthquake has a magnitude of 6.0 Mw, and the location is 14.59 km from the epicenter, indicating strong shaking. The VS30 value of 354 m/s suggests moderate amplification. The area consists of general buildings, and the visual context shows industrial/commercial structures that appear relatively robust. High population density (9,318.61 people per sq km) increases risk, but affluence and education suggest better infrastructure. Considering these factors, the damage is expected to be noticeable but not severe.

Qwen-2.5-32B | 2014 Napa | Predicted MMI: VIII

The earthquake has a magnitude of 6.0 Mw and occurred 10.24 km from the epicenter. The firm ground conditions (VS30 = 508 m/s) reduce shaking intensity. The area is sparsely populated (23.04 people per sq km), with low urbanization (2.97 %) and high affluence (\$93,750 median income). The visual context shows a rural landscape with minimal infrastructure, indicating low vulnerability to structural damage. Based on these factors, the expected damage level is moderate.

Table 5: Samples of LLM outputs with predicted MMI reasoning for selected earthquake scenarios.