

# BuildSafe: A Generative AI Framework for Joint Construction Safety, Permit, and Community-Impact Reasoning from Linked NYC DOB-OSHA-311 Data

1 **Shahid Ali<sup>1\*</sup>, Ammar Alzarrad<sup>2</sup>, Hwapyeong Song<sup>3</sup>, Husnu S. Narman<sup>4</sup>**

2 <sup>1</sup>Department of Computer Science, Marshall University, One John Marshall Drive, Huntington,  
3 WV 25755; e-mail: [shahidali@marshall.edu](mailto:shahidali@marshall.edu)

4 <sup>2</sup>Department of Civil Engineering, Marshall University, One John Marshall Drive, Huntington, WV  
5 25755; E-mail: [alzarrad@marshall.edu](mailto:alzarrad@marshall.edu)

6 <sup>3</sup>Department of Computer Science, Marshall University, One John Marshall Drive, Huntington,  
7 WV 25755; e-mail: [song24@marshall.edu](mailto:song24@marshall.edu)

8 <sup>4</sup>Department of Computer Science, Marshall University, One John Marshall Drive, Huntington,  
9 WV 25755; e-mail: [narman@marshall.edu](mailto:narman@marshall.edu)

10 **\* Correspondence:**

11 Shahid Ali  
12 [shahidali@marshall.edu](mailto:shahidali@marshall.edu)

13 **Keywords: Generative Artificial Intelligence; Large Language Models; Construction Safety;**  
14 **Urban Governance; Building Permits; Risk Prediction; New York City.**

15 **Abstract**

16 BuildSafe investigates whether domain-adapted generative models can turn routine administrative text  
17 permit filings, incident narratives, and citizen complaints into early-warning signals for construction  
18 risk in New York City. The study links Department of Buildings (DOB) records, OSHA enforcement  
19 cases, and 311 complaints into a shared schema. It uses Gemini 1.5 Flash to generate three-section  
20 annotations covering job-site hazards, permit/code requirements, and community impacts. Three open  
21 models (Gemma-3-1B, Llama-3.2-3B, and Mistral-7B-Instruct-v0.3, 4-bit) are then fine-tuned with  
22 LoRA on this corpus using a unified prompt template, and evaluated with BERTScore, ROUGE,  
23 BLEU, METEOR, and bootstrap confidence intervals on a 2,833-record held-out test set, and  
24 supplemented with a 150-record two-expert consensus validation subset used for targeted human  
25 evaluation of factual accuracy, completeness, and hallucinations. Across the short-run LoRA  
26 adaptation setting, all models showed small train-validation-test gaps, suggesting stable proof-of-  
27 concept adaptation without obvious overfitting rather than evidence of full convergence. Llama-3.2-3B  
28 is the clear leader on lexical overlap metrics, Mistral-7B-Instruct delivers the strongest semantic  
29 fidelity and most structured narrative outputs, and Gemma-3-1B offers competitive semantic  
30 performance at a fraction of the compute cost, suggesting potential value for edge or real-time triage  
31 workloads. A qualitative case study of Manhattan permit records illustrates how these models produce  
32 multi-section safety narratives and reveals trade-offs among completeness, specificity, and verbosity.  
33 The paper highlights key limitations, most notably the reliance on machine-generated reference labels,  
34 the limited scale of human expert evaluation, and the NYC-specific regulatory context, and proposes  
35 a roadmap for future operational validation rather than claiming demonstrated deployment impact.

36 Priorities include richer and better-linked data sources, retrieval-augmented and cascaded inference  
37 architectures, impact- and fairness-oriented evaluation in agency pilots, and human-in-the-loop  
38 governance suited to high-stakes public-sector AI systems.

## 39 **1 Introduction**

40 The construction sector remains one of the most hazardous industries worldwide, with persistent  
41 challenges around accident prevention, regulatory compliance, and community disruption (Tang, 2024;  
42 Usama et al., 2024). In dense urban environments such as New York City (NYC), these challenges  
43 intensify: projects must navigate complex building code and permitting requirements while minimizing  
44 impacts on surrounding neighborhoods (Awolusi et al., 2022; Ceccato, 2013). Empirical analyses of  
45 construction injuries and incident patterns show that reactive safety management based on post hoc  
46 investigation and fragmented data has not delivered sustained improvements in safety performance  
47 (Awolusi et al., 2022; Pilskog Orvik, 2024; Tixier and Hallowell, 2023).

48 Cities like NYC now publish rich, complementary data describing construction activity and its  
49 consequences. NYC Department of Buildings (DOB) housing and permitting datasets provide free-text  
50 job descriptions and structured attributes for new buildings, major alterations, and demolitions (New  
51 York City Department of City Planning, n.d.; Nycdb, n.d.). OSHA enforcement and injury-reporting  
52 data add incident narratives, cited standards, and employer-reported injuries and illnesses  
53 (Occupational Safety and Health Administration, 2025; Open Knowledge Foundation, n.d.). NYC 311  
54 service requests capture resident complaints about noise, unsafe building conditions, street  
55 obstructions, and related quality-of-life issues in fine spatial and temporal detail (City of New York,  
56 2026; Mulligan et al., 2019; Tussey and Yan, 2025). Individually, these datasets support focused  
57 analyses of incidents, delays, or community response; collectively, they describe a lifecycle from  
58 permit filing to realized safety outcomes and neighborhood complaints (Ceccato, 2013; Mulligan et  
59 al., 2019; Zou and Ergan, 2019). However, they are maintained in separate systems with heterogeneous  
60 schemas and are rarely linked at scale, limiting the use of historical patterns when triaging new permit  
61 applications (Gao et al., 2023; Kamil et al., 2024; Yang et al., 2017).

62 Recent AI research in construction safety has advanced along three main streams: vision and sensor-  
63 based monitoring of unsafe site conditions, statistical or machine-learning prediction of accidents and  
64 violations, and emerging generative AI systems that summarize or explain safety risks. These  
65 approaches have improved task-specific decision support, but they remain limited as foundations for  
66 pre-construction urban governance. Vision and sensor systems typically operate after work has begun  
67 and depend on observable site conditions. Traditional predictive models often estimate isolated  
68 outcomes such as injury likelihood, violation risk, or complaint frequency. Early LLM-based systems  
69 can generate safety guidance, but they rarely connect permit descriptions, enforcement histories, and  
70 neighborhood complaint signals within a unified governance workflow. Thus, the main limitation is  
71 not technical immaturity but misalignment between model design and the cross-agency information  
72 environment of public construction oversight. BuildSafe addresses this limitation by shifting the unit  
73 of analysis from isolated safety events to linked permit-enforcement-complaint records that support  
74 earlier and more explainable risk triage.

75 BuildSafe differs from prior AI-based construction safety frameworks in four specific ways. First, it  
76 treats construction safety as a linked urban-governance problem rather than a single-task hazard  
77 classification or accident prediction problem. Second, it integrates DOB permit records, OSHA  
78 enforcement and injury narratives, and 311 complaint histories into one administrative-civic corpus,  
79 whereas most prior systems analyze safety, permitting, or community disruption separately. Third, it

80 formulates the task as tri-section generative reasoning: a single prompt-output template jointly  
81 produces job-site hazard predictions, likely permit/code requirements, and potential community-impact  
82 signals. Fourth, it evaluates multiple open-weight LLMs under a common LoRA/QLoRA fine-tuning  
83 protocol with held-out testing, bootstrap confidence intervals, ablation analysis, zero-shot comparison,  
84 and targeted human validation. This positioning shifts the contribution from another isolated safety-  
85 prediction model toward an integrated, city-scale framework for explainable construction-risk  
86 governance.

87 This framing also connects BuildSafe to computational research on urban resilience. Maksoud et al.  
88 (2024a) show that computational design and multi-objective optimization can support flood-resilient  
89 urban habitats by combining performance simulation, environmental adaptation, and data-driven  
90 design evaluation. BuildSafe extends this broader logic from environmental design to construction-  
91 safety governance by using linked administrative and civic data to anticipate site-level hazards and  
92 neighborhood-scale impacts before risks fully materialize. This study therefore focuses on NYC  
93 construction job filings and asks:

- 94 • **RQ1.** Can a single generative model, fine-tuned on linked DOB-OSHA-311 data, generate  
95 useful joint predictions of hazards, permitting needs, and community impacts directly from  
96 DOB job descriptions?
- 97 • **RQ2.** Does a unified generative pipeline trained on tri-task prompts produce internally coherent  
98 outputs across the hazard, permit, and community-impact sections, and does the shared prompt-  
99 output template maintain semantic consistency across sections?

100 The working hypothesis is that a unified generative pipeline, trained on prompts that bundle DOB job  
101 attributes with spatio-temporal context from OSHA incidents and nearby 311 complaint histories, can  
102 learn cross-task regularities that are difficult to capture when each task is modeled separately (Gao et  
103 al., 2023; Tussey and Yan, 2025; Zou and Ergan, 2019). Examples include recurrent combinations of  
104 project scope and location that simultaneously elevate fall or struck-by risk, trigger specific NYC  
105 Building Code sections and permit types, and correlate with after-hours noise or dust complaints. By  
106 using a shared prompt-output template during fine-tuning, the model can exploit common lexical and  
107 semantic structure work types, construction methods, regulatory references, and neighborhood  
108 characteristics across the hazards, permits, and community-impacts sections, potentially improving  
109 both semantic fidelity and internal coherence (Mohamed et al., 2025; Rasheed et al., 2024; Yoo et al.,  
110 2024).

111 BuildSafe is positioned as an experimental framework and benchmark study rather than a deployed  
112 production system. The paper makes three contributions: (1) constructs a 90,000-prompt corpus by  
113 spatially and temporally linking NYC DOB housing jobs, OSHA enforcement and injury-reporting  
114 data, and NYC 311 service requests, and by generating structured three-section annotations (hazards,  
115 permits, community impacts) with a consistent template suitable for supervised fine-tuning (City of  
116 New York, 2026; New York City Department of City Planning, n.d.; Occupational Safety and Health  
117 Administration, 2025); (2) fine-tunes three open-weight LLMs Gemma-3-1B, Llama-3.2-3B, and  
118 Mistral-7B-Instruct (4-bit) using parameter-efficient strategies (LoRA/QLoRA) via Unsloth, creating  
119 three BuildSafe variants that differ in capacity and deployment profile (Farabet, 2025; Mistral AI,  
120 2024; Unslothai, 2025b); and (3) evaluates these models on a held-out DOB test set using standard  
121 text-generation metrics (BERTScore, ROUGE, BLEU, METEOR), analyzing complementary  
122 strengths in semantic similarity, lexical precision, and computational efficiency in this tri-task  
123 generative setting (Jiao, 2024; Mohamed et al., 2025; Wan et al., 2024). Beyond traditional n-gram  
124 metrics such as BLEU and ROUGE, embedding-based measures like BERTScore have been shown to

125 correlate more strongly with human judgments on generative language tasks (Zhang et al., 2020).  
126 Finally, the paper discusses how this framework could support municipal workflows, for example by  
127 flagging high-risk projects for prioritized review, assisting permit desk staff with code citations, and  
128 highlighting projects likely to generate specific community complaints while drawing on recent work  
129 in responsible and trustworthy AI to identify key governance requirements, including safety-aware  
130 fine-tuning, bias and fairness auditing, and appropriate regulatory oversight (Agapiou, 2024;  
131 Behzadan, 2024; Hsu et al., 2024).

## 132 **2 Literature Review**

133 The literature review is organized around the progression from traditional construction-safety risk  
134 assessment to data-driven prediction, generative AI, geospatial urban analytics, and trustworthy AI  
135 governance. Rather than treating these areas as separate bodies of work, the review emphasizes how  
136 each contributes only part of the BuildSafe problem: early identification of hazards, interpretation of  
137 permit and code requirements, and anticipation of community impacts from linked urban data.

### 138 **2.1 Construction Safety - Scale and Persistence**

139 Construction remains one of the deadliest industries in the United States, accounting for roughly one  
140 in five workplace fatalities and over 1,075 deaths in 2023, with the “Fatal Four” hazards (falls, struck-  
141 by, electrocutions, caught-in/between) responsible for nearly 60% of these deaths (Occupational Safety  
142 and Health Administration, 2023). Despite mature regulations and long-recognized hazard categories,  
143 fatality rates have not declined proportionally, underscoring the limits of traditional, reactive safety  
144 management. Expanding digital reporting has created large administrative datasets. OSHA’s Injury  
145 Tracking Application (ITA) mandates the electronic submission of Forms 300/300A/301 for high-  
146 hazard establishments, significantly increasing the structured nature of injury records, while sectoral  
147 efforts, such as the American Petroleum Institute’s guidelines, standardize occupational injury data  
148 (American Petroleum Institute, 2019; Occupational Safety and Health Administration, 2025).  
149 However, studies highlight that these records are rarely exploited for predictive modelling, and  
150 underreporting, fragmented subcontractor records, and project-based work further complicate  
151 longitudinal analysis (Piri and Panthi, 2024; Workers’ Compensation Trust, 2021). Evidence from  
152 incident analyses shows that leading indicators such as near-misses and violation patterns are present  
153 in existing data but not systematically leveraged (Awolusi et al., 2022), suggesting that the core  
154 problem is ineffective data utilization rather than data scarcity, an opportunity for AI and large  
155 language models (LLMs).

### 156 **2.2 From Probability - Impact Models to Data-Driven Risk Assessment**

157 Historically, construction risk assessment has been dominated by Probability-Impact (P-I) matrices and  
158 related expert-elicited methods, which yield ordinal risk scores, struggle with multiple objectives, and  
159 poorly capture interdependencies among risk factors (Taroun et al., 2011). Subsequent work introduced  
160 fuzzy logic, the Analytic Hierarchy Process, and Bayesian networks, yet systematic reviews still report  
161 methodological fragmentation and limited cross-study comparability (Kumi et al., 2024). Recent data-  
162 driven approaches address these limitations by fusing heterogeneous sources and learning latent  
163 relationships directly from data. Gao et al. (2023) showed that deep neural networks trained on multi-  
164 source quality data outperform traditional models, while Kamil et al. (2024) demonstrated analogous  
165 gains for heterogeneous safety data integration in the process industry. Mostofi and Toğan (2023)  
166 further increased representational power using multi-head attention graph neural networks, capturing  
167 relational accident structures but on relatively small datasets of uncertain generalizability. These  
168 limitations motivate BuildSafe’s move away from static, single-score risk assessment toward a linked,

169 data-driven generative framework that can reason jointly over hazards, permit requirements, and  
170 community impacts.

### 171 **2.3 AI and Machine Learning in Construction Safety Management**

172 AI and ML have increasingly been used to improve hazard detection, risk prediction, and decision  
173 support in construction safety. Reviews show progress in computer vision, NLP, predictive analytics,  
174 wearables, IoT sensors, drones, VR-based safety training, and AI-enabled compliance monitoring, but  
175 they also identify persistent barriers related to data quality, interpretability, standardization, and  
176 organizational readiness (Tang, 2024; Parekh and Mitchell, 2024; Savaş, 2025; Chandu et al., 2024;  
177 Hossain et al., 2025; Raliile and Haupt, 2020; HSI, 2025). Many proposed systems remain conceptual  
178 or proof-of-concept rather than field-validated, particularly in AI-driven risk management frameworks  
179 that are not yet integrated into routine agency or contractor workflows (Rasheed et al., 2024; Usama et  
180 al., 2024; Pilskog Orvik, 2024). This literature establishes the technical promise of AI in construction  
181 safety, while also showing that many systems remain disconnected from the administrative and civic  
182 data streams used in public construction governance. The remaining challenge is therefore not simply  
183 model development, but aligning AI methods with the cross-agency information environment in which  
184 public construction oversight actually occurs.

### 185 **2.4 Deep Learning and Predictive Modelling for Construction Accidents**

186 Deep learning expands construction-safety modeling by handling unstructured text, images, and sensor  
187 streams. GPT-based accident-prediction models and saliency methods have improved interpretability,  
188 while multi-company datasets show that broader shared data can outperform firm-specific models (Yoo  
189 et al., 2024; Tixier and Hallowell, 2023). Spatial and temporal context also improve predictive  
190 performance, with multi-order spatio-temporal models outperforming approaches that rely only on  
191 spatial or temporal features alone (Jiao, 2024). However, these studies typically focus on accident  
192 prediction or site-level monitoring rather than jointly modelling permitting requirements, safety  
193 enforcement signals, and neighborhood complaint patterns. BuildSafe extends this predictive literature  
194 by shifting the modeling target from isolated accident forecasting to joint reasoning across pre-  
195 construction permit language, historical enforcement evidence, and neighborhood complaint signals.

### 196 **2.5 Generative AI and Large Language Models in Construction**

197 Generative AI and LLMs shift the focus from pattern recognition to knowledge representation and  
198 explanation, which is critical for safety applications that require contextualized, actionable guidance  
199 rather than raw scores. Mohamed et al. (2025) performed a bibliometric analysis of generative AI in  
200 construction risk management, finding rapid growth since 2020 and categorizing benefits into  
201 technical, operational, technological, and integration domains, alongside nine risk categories spanning  
202 data quality, liability, and social acceptance. Wan et al. (2024) surveyed broader generative AI  
203 applications in the building industry (design, document generation, scenario simulation), but mostly at  
204 a conceptual level. Complementing this broader view, Maksoud et al. (2024b) show that image-  
205 generative AI tools such as Midjourney are already influencing architectural concept development,  
206 creative brainstorming, and design reasoning. However, their study also cautions that speed and  
207 creativity do not necessarily guarantee technical accuracy or input fidelity. A few studies propose  
208 working systems that combine vision and language. Hussain et al. (2026) built an intelligent assistant  
209 that integrates computer vision-based hazard detection with LLM-based safety advisories, while Tang  
210 and Luo (2025) showed that multimodal architectures can produce more nuanced recommendations  
211 than either vision or language alone.

212 Debate persists on trustworthiness. Behzadan (2024) argues that transparent architectures, explainable  
213 outputs, and user-centered design are prerequisites for adoption, and Mohamed et al. (2025) identify  
214 stakeholder trust as a key barrier. Proponents counter that LLMs need only outperform current manual  
215 practices, which are themselves error-prone and inconsistent (Yoo et al., 2024). BuildSafe contributes  
216 empirical evidence by fine-tuning domain-specific LLMs that can be evaluated against structured  
217 safety benchmarks. This motivates BuildSafe’s use of generative AI not as a generic advisory chatbot,  
218 but as a structured tri-section reasoning system whose outputs are constrained to hazards, permit/code  
219 requirements, and community impacts.

## 220 **2.6 Spatio-Temporal and Geospatial Approaches to Urban Construction Risk**

221 Urban construction risk is inherently spatial. Ceccato (2013) pioneered the use of GIS to uncover risk  
222 hotspots, corridors, and temporal clusters that non-spatial methods overlook. Dufitimana et al. (2025)  
223 combined ML with geospatial data to map socioeconomic vulnerability to natural hazards, and Yang  
224 et al. (2017) developed a georelational learning framework to integrate heterogeneous urban datasets,  
225 including permit and infrastructure records. Zou and Ergan (2019) showed that neighborhood-level  
226 features such as building density, complaint history, and land-use patterns predict construction-related  
227 disruptions, demonstrating that civic data infrastructures contain safety-relevant signals beyond site-  
228 level records. Together with Jiao’s (2024) findings on multi-order spatio-temporal models, this work  
229 motivates BuildSafe’s reliance on geographically resolved NYC data: spatial features are not ancillary  
230 metadata but core predictive variables.

## 231 **2.7 NYC Urban Data Sources for Construction Safety Research**

232 New York City offers an unusually rich open-data ecosystem for construction safety research,  
233 combining scale, transparency, and regulatory complexity. Key sources include NYC 311 service  
234 requests, DOB housing and permitting records, the NYCDB consolidated housing datasets, and U.S.  
235 Department of Labor enforcement data (Mulligan et al., 2019; City of New York, 2026; New York  
236 City Department of City Planning, n.d.; Nycdb, n.d.; Open Knowledge Foundation, n.d.). Zerkin (2006)  
237 documented the broader regulatory landscape of building performance in New York City, providing  
238 historical context for the permitting and inspection environment. 311 complaints provide  
239 neighborhood-level signals of disruption and potential unsafe conditions, DOB and housing records  
240 capture where and when construction occurs, and federal enforcement and ITA datasets provide  
241 standardized injury and violation data. However, these sources were not designed for joint analysis:  
242 they differ in geographic identifiers, temporal granularity, and categorical schemas. Prior work has  
243 begun to tackle geo-alignment and linkage challenges (Yang et al., 2017; Tussey and Yan, 2025), but  
244 a fully integrated, construction-safety-focused dataset has not been demonstrated. BuildSafe’s NYC  
245 DOB-OSHA-311 corpus directly addresses this gap. The importance of these datasets is therefore not  
246 only their scale, but their complementarity. DOB records describe proposed construction activity  
247 before work begins, OSHA records describe realized safety and compliance failures, and 311  
248 complaints capture neighborhood-level disruption that may not appear in formal safety records. Prior  
249 urban analytics studies have shown the value of geospatial and civic data, but they rarely convert these  
250 linked administrative traces into a structured generative-reasoning task. BuildSafe uses this  
251 complementarity to position construction safety as a city-scale information-integration problem rather  
252 than a site-level monitoring problem alone. Section 4.7 later revisits this city-scale framing in relation  
253 to urban informatics, computational urban resilience, and AI-supported architectural design workflows.

## 254 **2.8 Transfer Learning and Cross-domain Adaptation**

255 Labelled construction safety data are scarce relative to the complexity of the tasks, making transfer  
256 learning (TL) particularly attractive. Junjia et al. (2024) classify TL techniques used in construction  
257 risk management as instance-based, mapping-based, network-based, and adversarial and document  
258 applications across computer vision, NLP, and expert systems, while highlighting ongoing challenges  
259 such as domain shift and the lack of benchmark datasets. Tixier and Hallowell’s (2023) multi-company  
260 findings support the broader TL insight that wider source knowledge improves target-domain  
261 performance. In LLMs, fine-tuning can be viewed as network-based transfer learning: a large model  
262 pre-trained on general text is adapted to construction safety language and reasoning using domain-  
263 specific data. BuildSafe operationalizes this by fine-tuning open-weight models on a curated corpus of  
264 NYC construction safety data. This directly supports BuildSafe’s use of LoRA/QLoRA as a resource-  
265 efficient domain-adaptation strategy for translating general-purpose open-weight LLMs into  
266 construction-safety reasoning models.

## 267 **2.9 Fine-Tuning Large Language Models for Domain-Specific Tasks**

268 Fine-tuning pre-trained LLMs has become the dominant strategy for domain specialization. Reviews  
269 by Weng (2024) and Lu et al. (2025) survey full fine-tuning, adapters, prompt tuning, and Low-Rank  
270 Adaptation (LoRA), emphasizing trade-offs between computational cost, adaptation depth, and  
271 catastrophic forgetting. Sagar (2025) shows that parameter-efficient methods like LoRA can deliver  
272 strong domain performance under tight resource constraints.

273 LoRA and its quantized variant QLoRA are especially attractive for academic settings. Zhao et al.’s  
274 (2024) LoRA Land report demonstrates that LoRA-adapted models can rival GPT-4 on specialized  
275 benchmarks, and frameworks such as Unsloth provide substantial speed and memory savings  
276 (Unslothai, 2025b; Unsloth, 2026), enabling fine-tuning of multi-billion-parameter models on single  
277 GPUs. This ecosystem underpins BuildSafe’s resource-efficient fine-tuning strategy. Safety  
278 preservation during fine-tuning is a critical concern in safety-critical domains. Choi et al. (2024)  
279 propose safety-aware fine-tuning methods that maintain alignment, and Hsu et al. (2024) introduce  
280 Safe LoRA, constraining weight updates to a “safe” subspace. While BuildSafe does not yet implement  
281 these methods, their principles inform conservative hyperparameter choices and structured prompt  
282 templates, and they define a clear path for future safety-aware domain adaptation.

## 283 **2.10 Base Model Selection - Open-Weight Architectures**

284 Among open-weight models, Mistral-7B and Google’s Gemma families are particularly relevant for  
285 resource-constrained research. Mistral-7B combines strong instruction-following performance with  
286 architectural innovations such as grouped-query and sliding-window attention, and is widely available  
287 in 4-bit quantized form suitable for single-GPU fine-tuning (Mistral AI, 2024; Data Science Dojo,  
288 2025; Hugging Face, 2025a; Miller, 2024; Ubai, 2024). Gemma models offer competitive  
289 performance and tight integration with Google’s tooling, with versions explicitly designed to run on  
290 single GPUs or TPUs (Team G. et al., 2024; Farabet, 2025). For BuildSafe’s core task, structured  
291 extraction and generation of risk-relevant information from DOB job descriptions and related context,  
292 these architectures provide a practical balance between capacity, performance, and deployability, while  
293 illustrating that domain-adapted LLMs no longer require proprietary, closed-weight bases.

## 294 **2.11 Regulatory Frameworks, Ethics, and Trustworthy AI in Construction**

295 Deploying AI for safety-critical decisions raises socio-legal and ethical challenges. Agapiou (2024)  
296 reviews responsible AI in construction health and safety and concludes that existing legal frameworks  
297 do not yet adequately address algorithmic bias, opacity, and liability allocation when AI contributes to

298 safety failures. More broadly, Huang et al. (2024), Sargiotis (2024), and Wörsdörfer (2025) highlight  
299 the rapid expansion of AI applications in civil engineering and the uneven global trajectories of AI  
300 regulation, reinforcing the need for stronger governance in safety-critical built-environment systems.  
301 Yang et al. (2024) propose a lifecycle framework for trustworthy AI in construction, spanning  
302 planning, data collection, algorithm development, deployment, and maintenance, and mapping  
303 principles such as safety, explainability, transparency, and fairness to concrete actions. For fine-tuned  
304 LLMs like BuildSafe, this implies that predictive accuracy must be accompanied by transparency about  
305 training data and limitations, mechanisms for human oversight, and alignment with data privacy and  
306 liability norms. Accordingly, the present study treats ethics and trustworthiness not only as background  
307 concerns but as deployment requirements, with Section 4.11 translating these principles into  
308 measurable fairness audits, privacy and cybersecurity controls, audit logging, and accountability  
309 procedures.

## 310 **2.12 Research Gaps and Rationale for the Present Study**

311 Taken together, the literature shows that construction-safety AI has progressed from expert-driven risk  
312 matrices to machine-learning prediction, computer-vision monitoring, geospatial analysis, and early  
313 generative-AI decision support. However, these streams remain only partially connected. Sensor,  
314 wearable, drone, and vision-based systems are most useful after work has begun, while permit-review  
315 agencies require risk signals earlier, when administrative records, project descriptions, location  
316 attributes, enforcement histories, and complaint patterns are the main available evidence. Traditional  
317 ML and deep learning studies improve prediction, but they usually treat accident occurrence, violation  
318 risk, and complaint frequency as separate targets rather than as interrelated elements of urban  
319 construction governance. Generative-AI studies broaden the scope by producing explanatory outputs,  
320 but many remain conceptual, rely on generic prompting, or do not integrate permitting, enforcement,  
321 and civic-impact data into a reproducible training and evaluation workflow.

322 This synthesis leaves four focused gaps that motivate BuildSafe. First, few studies evaluate working,  
323 domain-adapted LLMs against a structured construction-safety benchmark. Second, DOB permits,  
324 OSHA records, 311 complaints, and geospatial indicators are usually analyzed separately rather than  
325 as a linked training corpus. Third, resource-efficient LLM adaptation methods such as LoRA and  
326 QLoRA remain underexplored for construction-safety reasoning, despite their relevance for academic  
327 and public-sector settings with limited computing resources. Fourth, safety-aware adaptation and  
328 governance-oriented evaluation remain insufficiently operationalized, even though unsupported AI  
329 outputs could affect permitting, inspection, and public-sector decision workflows. BuildSafe addresses  
330 these gaps by testing whether linked administrative, enforcement, and civic data can support a unified  
331 generative framework for construction-safety reasoning.

## 332 **3 Methodology**

### 333 **3.1 Data Acquisition and Sources**

#### 334 **3.1.1 NYC DOB Housing Database**

335 The core inventory of construction activity is the New York City Department of City Planning's  
336 Housing Database, which consolidates Department of Buildings (DOB) job filings affecting residential  
337 unit counts across the five boroughs beginning 1 January 2010 (New York City Department of City  
338 Planning, n.d.). The database distinguishes three primary job types: new buildings, Alteration Type 1  
339 (major alterations that change dwelling-unit counts), and demolitions, and provides project-level  
340 records with attributes such as borough, block and lot, address, job type, and pre- and post-development

341 unit counts, all geocoded to support spatial joins. The Housing Database is distributed via NYC Open  
342 Data and BYTES of the BIG APPLE as CSV files. It is mirrored in a version-controlled PostgreSQL  
343 schema by the open-source nycdb project, which simplifies automated ingestion (Nycdb, n.d.). In  
344 BuildSafe, this database serves as the authoritative list of construction projects against which other  
345 safety-relevant signals are spatially and temporally aligned.

### 346 **3.1.2 OSHA Enforcement and Injury-Tracking Data**

347 Workplace safety outcomes are characterized using enforcement and injury data released by the U.S.  
348 Department of Labor. The US DOL enforcement data portal exposes OSHA inspection, violation, and  
349 penalty records as bulk CSV exports and APIs, linking inspection identifiers, cited standards, employer  
350 information, and penalties (Open Knowledge Foundation, n.d.). OSHA’s Injury Tracking Application  
351 (ITA) complements these tables with establishment-specific illness and injury records from Forms 300,  
352 300A, and 301 (Occupational Safety and Health Administration, 2025). Narrative fields in these  
353 datasets, such as alleged violation descriptions and incident summaries, encode fine-grained  
354 information about hazard types and failure modes that structured variables alone cannot capture. These  
355 narratives are used to expose the model to authentic descriptions of construction-related incidents and  
356 regulatory framing.

### 357 **3.1.3 NYC 311 Service Requests**

358 To capture community-level signals of construction impact, the study incorporates NYC 311 service  
359 requests, the city’s centralized non-emergency reporting system. The dataset contains millions of  
360 complaints annually about noise, dust, unsafe building conditions, and street obstructions, with each  
361 record including timestamps, complaint type, responsible agency, and geocoded location (Mulligan et  
362 al., 2019; City of New York, 2026). The data are regularly updated on NYC Open Data and support  
363 both bulk downloads and API access, enabling large-scale spatio-temporal analysis. Complaint  
364 categories related to buildings and noise are particularly relevant for understanding public perceptions  
365 of construction activity. By aggregating complaint counts within spatial buffers around DOB project  
366 locations, BuildSafe infers patterns of community disruption associated with different project types.

## 367 **3.2 Data Preprocessing and Integrations**

### 368 **3.2.1 Ingestion and Encoding Normalization**

369 Fifteen heterogeneous CSV files spanning NYC 311, DOB housing records, OSHA enforcement  
370 tables, and auxiliary datasets constitute the raw input corpus. Each file is ingested into pandas using a  
371 two-stage decoding strategy: the pipeline first attempts UTF-8 and, on failure, retries with Latin-1  
372 while skipping corrupted lines. For each file, the script logs the encoding used, row counts, and any  
373 errors, providing traceability for later data-quality audits.

### 374 **3.2.2 Schema Harmonization**

375 Because source systems use divergent naming conventions, a predefined mapping aligns each file’s  
376 columns to a common minimal schema. This mapping harmonizes semantically equivalent fields (for  
377 example, boro and borough, job\_desc and job\_description) and drops attributes not needed for risk  
378 reasoning. Records lacking critical fields such as geographic identifiers or core descriptive text are  
379 removed, and remaining nulls in retained columns are converted to the literal string “N/A” to ensure  
380 syntactic completeness when forming prompts. The result is a consistent schema across all input tables  
381 with sufficient expressive power for safety assessment.

### 382 **3.2.3 Sampling Strategy and Full-Corpus Construction**

383 An initial deterministic sampling step selects up to five records per file as a quality-assurance subset  
384 to validate the annotation prompt template and verify response quality from Gemini 1.5 Flash. Using  
385 fixed random seeds, this per-file cap yields a balanced sample across datasets and complaint or incident  
386 types. After confirming that the three-section annotation structure (hazards, permits, community  
387 impacts) was reliably produced, the pipeline was extended to all records across the 15 source files that  
388 passed schema harmonization and null-filtering. In the full annotation run, every eligible record was  
389 submitted to Gemini 1.5 Flash using a fixed prompt template and a rate-limiting protocol. The resulting  
390 corpus comprises approximately 90,000 prompt-output pairs, each derived from a single DOB, OSHA,  
391 or 311 record's descriptive fields. The data were partitioned into training, validation, and test sets using  
392 a stratified random split with a fixed seed (3407), preserving the relative proportions of DOB, OSHA,  
393 and 311 sources across the splits. The test set contains 2,833 prompt-output pairs; the remaining records  
394 are divided between training and validation at roughly 80/20. Records with malformed or incomplete  
395 annotations, for example, outputs missing one or more of the three required sections, were excluded,  
396 with an overall exclusion rate below 2% of submitted records.

## 397 **3.3 Annotation Pipeline and Dataset Construction**

### 398 **3.3.1 Prompt Construction**

399 For each record, a helper routine constructs a multi-line text block in which each line follows the  
400 pattern field name: field value. When a dedicated description field exists (for example, a DOB job  
401 description or OSHA narrative), it is prominently included; otherwise, the composite block serves as  
402 the primary contextual description. This representation preserves column semantics while presenting  
403 the model with a human-readable summary. The text block is embedded into a fixed prompt template  
404 instructing the model to: (1) assess the overall risk level associated with the described situation or  
405 project; (2) identify relevant permits, regulations, or standards likely to apply; (3) anticipate potential  
406 worker safety and community impacts; and (4) flag obvious compliance concerns or missing  
407 safeguards. This standardized design ensures that annotations across heterogeneous records are  
408 comparable and suitable as labels for supervised fine-tuning.

### 409 **3.3.2 LLM Invocation and Rate Limiting**

410 The annotation pipeline calls Google's Gemini 1.5-Flash model via a configured client, which sends  
411 one prompt per record. To comply with usage policies (e.g., a maximum of about 15 requests per  
412 minute), the script enforces a fixed delay of approximately 4 seconds between calls. Structured  
413 exception handling logs transport errors, rate-limit responses, and other API faults; failed records are  
414 retried with backoff or marked for exclusion depending on error type.

### 415 **3.3.3 Result Aggregation and Final Dataset**

416 Each successful annotation is stored with its provenance (source file, row index, and exact prompt  
417 text). After processing all records, the annotations are consolidated into an intermediate CSV file that  
418 includes both raw prompts and model-generated outputs. A final projection step extracts two canonical  
419 columns for fine-tuning: prompt and output. Where available, prompt is set to the original human-  
420 written description; otherwise, it defaults to the composite text block. The output column stores the  
421 unmodified annotation string returned by Gemini. The resulting two-column dataset is model-agnostic  
422 and suitable for supervised fine-tuning of any compatible LLM.

### 423 3.3.4 Annotation Reliability and Silver-Standard Supervision

424 The Gemini-generated annotations are treated in this study as silver-standard supervision rather than  
425 as definitive human ground truth. This distinction is central to the experimental design because the  
426 annotation task requires contextual safety reasoning, regulatory interpretation, and community-impact  
427 inference, all of which may be affected by model hallucinations, omissions, or overgeneralizations.  
428 Gemini 1.5 Flash was used because it provided a scalable, consistent mechanism for generating  
429 structured three-section annotations across heterogeneous DOB, OSHA, and 311 records using a fixed  
430 prompt template. The annotation task was constrained rather than open-ended: each administrative  
431 record was mapped into three predefined sections: job-site hazards, permit/code requirements, and  
432 community impacts, so the model functioned primarily as a structured extraction and safety-reasoning  
433 annotator. This use is consistent with recent construction and built-environment literature that treats  
434 generative AI as a tool for document interpretation, compliance-oriented reasoning, and risk-  
435 management support, while emphasizing the need for safeguards against hallucination, bias, and  
436 overgeneralization (Mohamed et al., 2025; Tan et al., 2024; Wan et al., 2024). However, the resulting  
437 labels should be interpreted as machine-assisted reference outputs suitable for proof-of-concept  
438 supervised fine-tuning, not as fully verified expert regulatory determinations.

439 The reliability of the silver-standard corpus was strengthened through four safeguards. First, all records  
440 were annotated using a fixed prompt template requiring the same three sections: hazards, permit/code  
441 requirements, and community impacts. This reduced variation in output structure and allowed the  
442 downstream models to learn a consistent tri-task response format. Second, each annotation was stored  
443 with source-file provenance, row index, and the exact prompt text, allowing every generated label to  
444 be traced back to the original administrative record. Third, malformed annotations were excluded when  
445 they lacked one or more required sections, resulting in an exclusion rate of less than 2% of submitted  
446 records. Fourth, the automated evaluation was supplemented by a stratified 150-record expert subset,  
447 independently annotated by two domain experts. For this subset, A.A., a civil engineering faculty  
448 member with construction safety and permitting expertise, and H.S.N., a computer science researcher  
449 with expertise in AI systems, independently authored three-section safety narratives covering hazards,  
450 permit/code requirements, and community impacts using the same output template as BuildSafe. Both  
451 experts were blinded to the Gemini-generated teacher labels and to the fine-tuned model outputs during  
452 annotation. Disagreements were resolved through discussion, and the resulting consensus narratives  
453 were treated as a human-reference benchmark for the 150-record validation subset.

454 Accordingly, the automated metrics reported on the full held-out test set should be interpreted as  
455 measuring agreement with a consistent teacher model across a large corpus, not as direct proof of  
456 absolute regulatory correctness. BERTScore, ROUGE, BLEU, and METEOR, therefore, quantify  
457 whether the fine-tuned models reproduce the structured safety-reasoning pattern encoded in the silver-  
458 standard references. They do not, by themselves, establish that every hazard, code requirement, or  
459 community-impact statement is legally or operationally correct. However, the 150-record two-expert  
460 consensus subset provides a separate human-reference benchmark, not derived from Gemini outputs,  
461 for evaluating whether model outputs remain factually accurate, complete, and free from unsupported  
462 regulatory or hazard claims. This design allows the study to combine scalable silver-standard  
463 supervision for model adaptation with a smaller but stronger human-grounded validation set. The  
464 purpose of the Gemini-generated labels is to test whether open-weight LLMs can be aligned with a  
465 consistent construction-safety reasoning format using linked DOB-OSHA-311 data. At the same time,  
466 the expert consensus subset provides a more rigorous check on domain plausibility and hallucination  
467 behavior. Any operational deployment would still require human review, retrieval from current  
468 regulatory sources, larger external expert validation, reporting of inter-rater agreement before

469 consensus resolution, and post-deployment monitoring before outputs could be used in safety-critical  
470 public-sector workflows.

471

## 472 **3.4 Model Selection and Experimental Setup**

### 473 **3.4.1 Model Family Rationale**

474 Three open-weight model families were selected for BuildSafe: Gemma, Llama, and Mistral, spanning  
475 a spectrum from lightweight to mid-sized architectures. Gemma models, derived from Gemini  
476 research, are designed for efficient deployment on single GPUs or TPUs while maintaining strong  
477 language understanding, making smaller variants attractive for low-latency triage (Team G. et al.,  
478 2024; Farabet, 2025). Llama 3.2 models provide mid-scale capacity (1B-3B parameters) and strong  
479 instruction-following performance suitable for more nuanced reasoning on modest hardware  
480 (Grattafiori et al., 2024).

481 Mistral-7B-Instruct-v0.3, a 7.3-billion-parameter decoder-only transformer, has been shown to  
482 perform strongly on a range of benchmarks and incorporates grouped-query attention and sliding-  
483 window attention for efficient long-sequence processing (Mistral AI, 2024; Data Science Dojo, 2025).  
484 The availability of a 4-bit quantized variant via bitsandbytes and Unsloth further reduces hardware  
485 requirements (Hugging Face, 2025a; Unslothai, 2025b). In the experiments, Gemma-3-1B serves as  
486 the most compact generative model for fast hazard and permit triage, Llama-3.2-3B as a mid-tier  
487 reasoning model, and Mistral-7B-Instruct-4bit as the primary candidate for high-fidelity risk  
488 assessment and narrative generation.

### 489 **3.4.2 Parameter-Efficient Fine-Tuning with Unsloth**

490 Adapting these base models to the construction safety domain relies on parameter-efficient fine-tuning  
491 (PEFT) via the Unsloth framework. Unsloth integrates Low-Rank Adaptation (LoRA) and quantized  
492 LoRA (QLoRA) to reduce VRAM usage and training time, with reported speedups and memory  
493 savings over naive fine-tuning (Hugging Face, 2025b; Unslothai, 2025a). This aligns with broader  
494 findings that LoRA-style adaptations can deliver domain-specific performance comparable to full fine-  
495 tuning at a fraction of the computational cost (Zhao et al., 2024; Sagar, 2025; Weng, 2024; Lu et al.,  
496 2025). In this study, LoRA modules are applied to the language, attention, and feed-forward layers  
497 while keeping the vast majority of base weights frozen, offering a practical compromise between  
498 adaptation depth and computational feasibility for academic-grade hardware.

### 499 **3.4.3 Fine-Tuning Configuration**

500 Supervised fine-tuning is conducted using the two-column prompt-output dataset, along with the  
501 SFTTrainer implementation provided by Unsloth. The configuration employs LoRA with a rank of 8  
502 and alpha of 8, without dropout or additional bias parameters, and is applied exclusively to text-based  
503 modules. Training uses a per-device batch size of 2 with gradient accumulation set to 4, resulting in an  
504 effective batch size of 8. Optimization is performed using the 8-bit AdamW optimizer with a learning  
505 rate of  $2 \times 10^{-4}$  for short exploratory runs, while a reduced rate of  $2 \times 10^{-5}$  is recommended for longer  
506 training schedules; weight decay is set to 0.01. The learning schedule used 5 warmup steps followed  
507 by linear decay across a maximum of 60 optimization steps. This short schedule should be interpreted  
508 as a constrained-compute, parameter-efficient adaptation test rather than as evidence of full  
509 convergence or globally optimized model performance. The purpose was to determine whether small  
510 LoRA updates could align pretrained open-weight models with the BuildSafe tri-task output schema

511 under academic-grade hardware constraints. Because the base models already encode broad linguistic  
512 and instruction-following capabilities, the adapters primarily learned the structured “hazards-permits-  
513 community impacts” response format and domain-specific lexical patterns rather than construction-  
514 safety reasoning from scratch. Metrics were logged at every training step, and train, validation, and test  
515 performance were compared as diagnostic checks for obvious instability or overfitting. However, these  
516 stability checks do not rule out the possibility that longer training, additional random seeds, alternative  
517 learning rates, checkpoint-based early stopping, or larger hyperparameter sweeps could further  
518 improve performance. Therefore, the reported results should be interpreted as reproducible proof-of-  
519 concept adaptation results, while a full convergence study remains necessary before claiming  
520 production-level optimization. A fixed random seed of 3407 was maintained throughout data loading,  
521 sampling, and training to support reproducibility. Accordingly, the reported metrics should be  
522 interpreted as evidence of short-run adapter alignment and output-format stabilization, not as evidence  
523 that the models have reached their maximum attainable domain performance.

### 524 **3.5 Human Evaluation Sub-Study**

525 Because the main automated metrics compare model outputs against Gemini 1.5 Flash-generated  
526 references on the full test set, a targeted human evaluation study was conducted to provide a separate  
527 human-reference benchmark, not derived from Gemini outputs, for a smaller subset of records. A  
528 stratified random sample of 150 records was drawn from the 2,833-record held-out test set, preserving  
529 the proportional distribution of source types (DOB, OSHA, 311) and job types (new building,  
530 Alteration Type 1, demolition). For this subset, two domain experts independently authored human  
531 reference narratives. A.A., a civil engineering faculty member with expertise in construction safety and  
532 building permitting, and H.S.N., a computer science researcher with expertise in AI systems, each  
533 produced three-section safety narratives covering Hazards, Permit/Code Requirements, and  
534 Community Impacts using the same template as BuildSafe’s outputs. Both experts were blinded to the  
535 Gemini-generated teacher labels and to the fine-tuned model outputs during this process.  
536 Disagreements between the two expert-authored narratives were resolved through discussion, and the  
537 final consensus narratives were treated as a human-reference benchmark for the 150-record subset.  
538 Because both experts are members of the author team, this consensus process strengthens reliability  
539 relative to a single-evaluator design but does not substitute for future external validation by  
540 independent agency reviewers, safety inspectors, or permitting professionals. Model outputs for this  
541 subset were assessed against the consensus human references using ordinal expert judgments for  
542 Factual Accuracy, Completeness, and Overall Usefulness, along with a hallucination assessment  
543 indicating whether the output referenced a non-existent or clearly inapplicable statute, code section,  
544 hazard category, or community-impact claim. Because this validation was designed as a targeted  
545 internal plausibility check rather than a fully powered external annotation study, the results are reported  
546 as directional expert-assessment trends rather than as definitive quantitative human-evaluation scores.

### 547 **3.6 Zero-Shot, Ablation, and Task-Structure Baselines**

548 To strengthen experimental rigor, three additional baseline and ablation analyses were included in the  
549 revised evaluation design. *First*, zero-shot baselines were evaluated by applying the original Gemma-  
550 3-1B, Llama-3.2-3B, and Mistral-7B-Instruct models to the held-out prompts without LoRA  
551 adaptation. These baselines test whether the observed BuildSafe performance reflects domain  
552 adaptation rather than only general instruction-following ability. *Second*, data-source ablations were  
553 conducted to examine the contribution of each linked urban data stream. Ablated variants were trained  
554 or evaluated using DOB-only, DOB+OSHA, DOB+311, and full DOB-OSHA-311 configurations.  
555 This design tests whether enforcement narratives and complaint histories provide distinct value beyond

556 permit descriptions alone. *Third*, task-structure ablations compared the unified tri-task model against  
557 single-task variants specialized for hazards, permit/code requirements, or community impacts. This  
558 comparison tests whether joint generation improves cross-section coherence or whether separate  
559 specialized models perform better on individual sections.

560 Traditional machine learning baselines such as logistic regression, support vector machines, random  
561 forests, and gradient-boosted trees were considered. However, they were not used as primary  
562 comparators because BuildSafe is not formulated as a single-label classification or scalar risk-scoring  
563 task. Its output is a structured, multi-section narrative that jointly explains hazard reasoning,  
564 permit/code interpretation, and community-impact assessment. A conventional ML baseline would  
565 therefore require converting the task into separate categorical or binary labels, such as a hazard  
566 category, a permit-escalation flag, or a high-complaint-risk indicator. Such a transformation would  
567 remove much of the generative reasoning and cross-sectional coherence that this study is designed to  
568 evaluate. Therefore, the most direct comparators for the present benchmark are zero-shot LLM outputs,  
569 data-source ablations, and task-structure ablations under the same prompt-output format. The study  
570 does not claim superiority over traditional ML methods on matched classification tasks; instead, future  
571 work should derive simplified classification or ranking targets from the BuildSafe corpus to enable fair  
572 comparison with logistic regression, support vector machines, random forests, gradient boosting, and  
573 retrieval-based non-generative systems.

### 574 **3.6.1 Baseline Implementation Details**

575 For the zero-shot comparison, the original unfine-tuned Gemma-3-1B, Llama-3.2-3B, and Mistral-7B-  
576 Instruct models were evaluated on the same held-out prompt set used for the fine-tuned BuildSafe  
577 models. The same input template, output-section requirements, evaluation scripts, and text-generation  
578 metrics were used so that differences reflected the effect of LoRA adaptation rather than differences in  
579 prompting or scoring. For data-source ablations, the DOB-only, DOB+OSHA, DOB+311, and full  
580 DOB-OSHA-311 configurations were compared using the same tri-task output structure to examine  
581 how each linked data stream affected hazards, permit/code reasoning, and community-impact  
582 narratives. For task-structure ablations, single-section variants were compared with the unified tri-task  
583 formulation to assess whether separating hazards, permits/code requirements, and community impacts  
584 improved section-specific accuracy or weakened cross-section coherence. These additional analyses  
585 test directional robustness and task-design sensitivity rather than final production optimization.

## 586 **4 Results and Discussion**

587 The results are interpreted in light of the annotation design. Because the full-scale automated metrics  
588 compare fine-tuned model outputs with Gemini-generated silver-standard references, they primarily  
589 measure consistency with a structured teacher model rather than verified regulatory correctness. The  
590 two-expert consensus subset provides a separate human-reference check on factual accuracy,  
591 completeness, usefulness, and hallucination behavior. Accordingly, the metric tables are treated as  
592 evidence of how well each model reproduces the BuildSafe tri-task reasoning format, not as proof that  
593 the outputs constitute legally validated regulatory determinations. Because BuildSafe generates open-  
594 ended three-section narratives, exact n-gram metrics such as BLEU, ROUGE, and METEOR are  
595 interpreted as indicators of lexical consistency. BERTScore is used as the primary semantic-similarity  
596 metric, while the remaining metrics help assess how closely each model follows the wording and  
597 structure of the silver-standard references. Detailed model-to-model interpretation is consolidated in  
598 Section 4.4, and qualitative output behavior is examined in Section 4.7.

### 599 **4.1 Gemma-3-1B**

600 Gemma-3-1B is the compact BuildSafe variant intended for rapid first-pass screening under limited  
 601 compute conditions. Table 1 shows that its train, validation, and test scores remain tightly aligned, with  
 602 test BERTScore F1 = 0.7703, indicating stable short-run adaptation without obvious overfitting. Figure  
 603 1 places Gemma within the broader model comparison, where its main value is low-compute semantic  
 604 screening rather than detailed regulatory phrasing.

605 **Table 1. Fine-Tuning Results of Gemma-3-1B**

Metric	Train	Validation	Test
BLEU	0.002507	0.002551	0.002444
BERTScore Precision	0.770643	0.770674	0.769623
BERTScore Recall	0.771463	0.771450	0.771280
BERTScore F1	0.770924	0.770929	0.770319
ROUGE-1	0.116625	0.116443	0.115389
ROUGE-2	0.048707	0.048967	0.047900
ROUGE-L	0.078756	0.078706	0.078486
ROUGE-LSum	0.112053	0.111912	0.111150
METEOR	0.020600	0.020636	0.020598

606  
 607 Gemma-3-1B’s main value is computational efficiency. Despite being the smallest model, it maintains  
 608 stable semantic performance across the train, validation, and test sets. Its weaker lexical-overlap scores  
 609 suggest that it is less appropriate for detailed regulatory phrasing or complete reviewer-facing  
 610 narratives, but useful for rapid first-pass screening where low latency and modest hardware  
 611 requirements are priorities.

612 **4.2 Llama-3.2-3B**

613 Llama-3.2-3B is the lexical precision leader among the BuildSafe variants. Table 2 shows that it  
 614 achieves the strongest surface-form alignment, including test ROUGE-1 = 0.1470 and METEOR =  
 615 0.0379, while maintaining high semantic similarity with test BERTScore F1 = 0.7742. This profile  
 616 makes Llama especially suitable for permit-review support, where consistent terminology and  
 617 standardized reviewer summaries are important.

618 **Table 2. Fine-Tuning Results of Llama-3.2-3B**

Metric	Train	Validation	Test
--------	-------	------------	------

BLEU	0.009324	0.009458	0.009143
BERTScore Precision	0.769804	0.769924	0.768811
BERTScore Recall	0.779874	0.780015	0.779851
BERTScore F1	0.774722	0.774850	0.774214
ROUGE-1	0.148643	0.149000	0.147028
ROUGE-2	0.054165	0.054723	0.053449
ROUGE-L	0.095700	0.095653	0.094863
ROUGE-LSum	0.139127	0.139381	0.137762
METEOR	0.037639	0.037821	0.037862

619 Llama-3.2-3B provides the strongest lexical consistency, with the highest BLEU, ROUGE-1, and  
620 METEOR scores among the three models. This profile makes it suitable for workflows where  
621 standardized wording, repeated permit-review language, and consistent reviewer summaries are  
622 important. Its strength is not that it produces the most elaborate explanations, but that it most closely  
623 follows the reference structure and terminology.

### 624 4.3 Mistral-7B-Instruct-v0.3 (4-bit)

625 Mistral-7B-Instruct provides the strongest semantic and narrative profile in the BuildSafe suite. Table  
626 3 shows the highest test BERTScore F1 value, 0.7747, despite weaker lexical-overlap scores than  
627 Llama. This pattern suggests that Mistral preserves meaning while using less reference-matched  
628 wording, making it more appropriate for complex or escalated cases that require coherent explanatory  
629 narratives rather than short standardized summaries.

630 **Table 3. Fine-Tuning Results of Mistral-7B-Instruct (4-bit)**

Metric	Train	Validation	Test
BLEU	0.00244	0.00247	0.00237
BERTScore Precision	0.7750	0.7751	0.7746
BERTScore Recall	0.7754	0.7755	0.7748
BERTScore F1	0.7752	0.7753	0.7747
ROUGE-1	0.1169	0.1168	0.1156
ROUGE-2	0.0492	0.0494	0.0483

ROUGE-L	0.0791	0.0790	0.0787
ROUGE-LSum	0.1124	0.1122	0.1114
METEOR	0.02045	0.02050	0.02046

631 Mistral-7B-Instruct provides the strongest semantic profile, achieving the highest BERTScore F1  
632 despite weaker lexical overlap than Llama. This pattern suggests that Mistral tends to preserve meaning  
633 while using more compressed or less reference-matched wording. It is therefore better suited to  
634 escalated cases where explanatory coherence matters more than exact reproduction of teacher-  
635 reference phrasing.

#### 636 4.4 Comparative Analysis

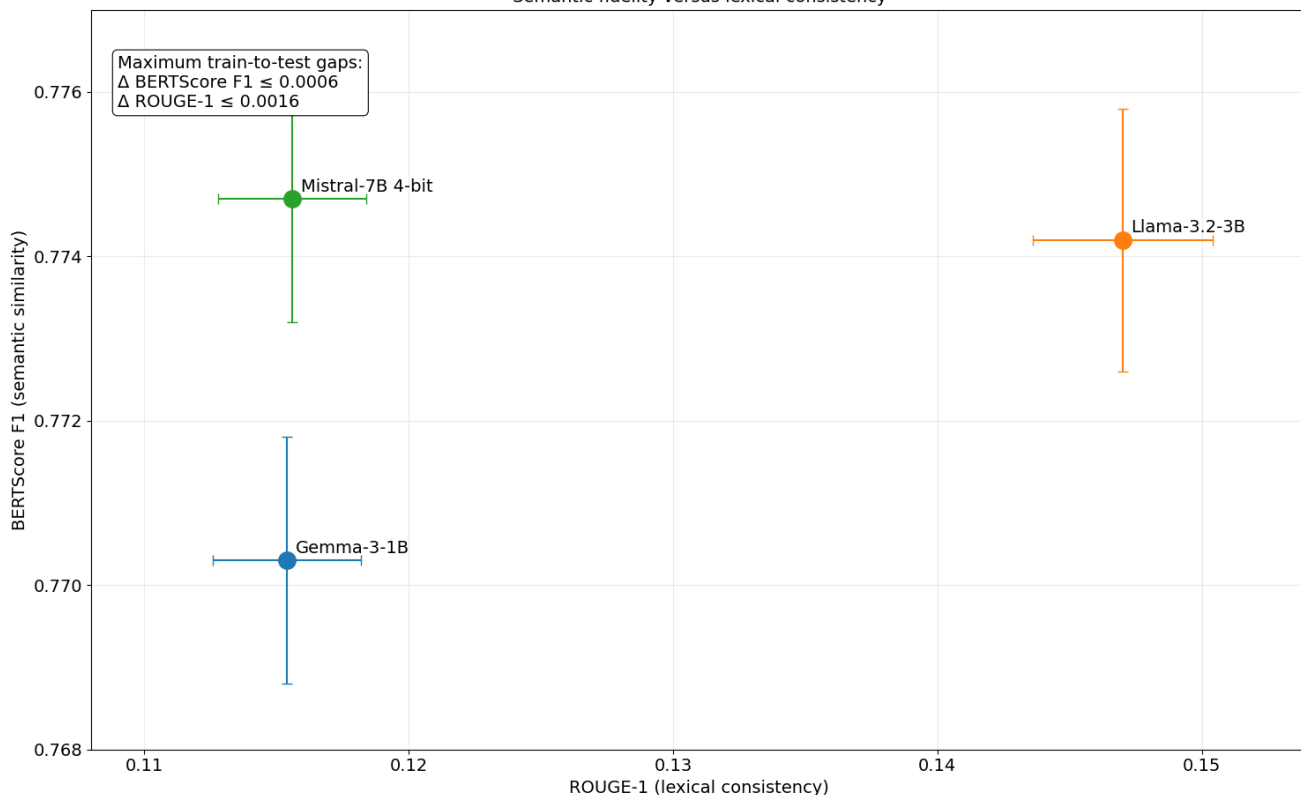
637 Rather than treating the three models as independent case descriptions, this section compares them  
638 along four dimensions that matter for construction-governance use: semantic fidelity, lexical  
639 consistency, generalization stability, and computational feasibility. Table 4 summarizes the main test-  
640 set metrics for Gemma-3-1B, Llama-3.2-3B, and Mistral-7B-Instruct, highlighting the trade-off  
641 between surface-form reproduction and meaning preservation.

642 **Table 4. Test-Set Metric Summary Across Models**

Model	BLEU	BERTScore F1	ROUGE-1	METEOR
Gemma-3 1B	0.002444	0.7703	0.1154	0.0206
Llama-3.2 3B	0.00914	0.7742	0.1470	0.0379
Mistral-7B 4-bit	0.00237	0.7747	0.1156	0.0205

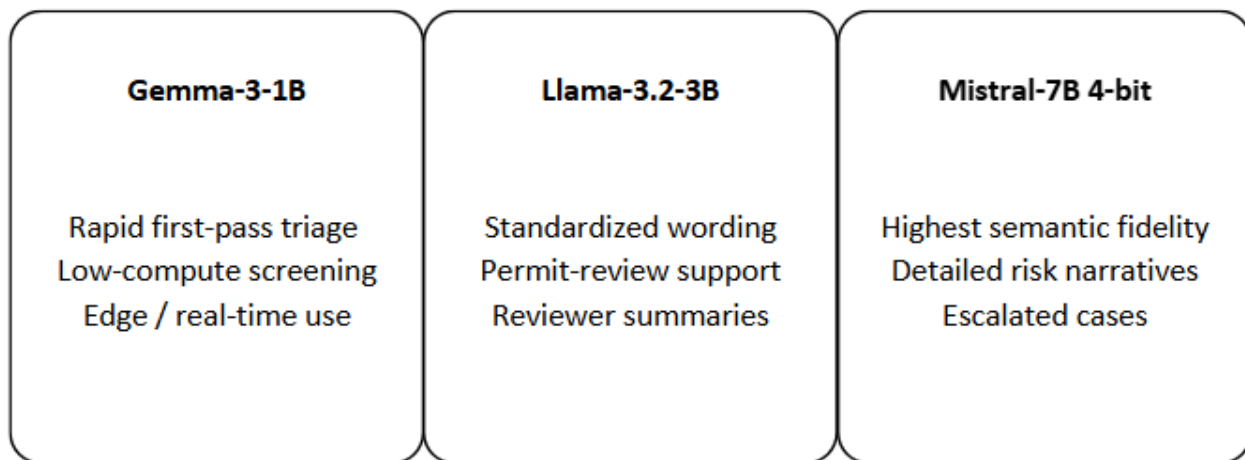
643

Semantic fidelity versus lexical consistency



644

Recommended Deployment Role



645

646 **Figure 1. Comparative performance and deployment positioning of BuildSafe model variants**

647 (A) BERTScore F1 plotted against ROUGE-1 for the three BuildSafe variants, with 95% bootstrap  
 648 confidence intervals. The plot highlights the trade-off between semantic fidelity and lexical  
 649 consistency: Llama-3.2-3B shows the strongest lexical overlap, while Mistral-7B-Instruct achieves  
 650 marginally higher semantic similarity. Models positioned toward the upper-right region of the plot  
 651 represent the strongest combined semantic and lexical balance; Llama-3.2-3B’s placement indicates  
 652 lexical dominance with strong semantic fidelity, whereas Mistral-7B-Instruct’s placement reflects  
 653 semantic depth with more compressed lexical overlap. (B) Recommended workflow roles implied by

654 the performance profiles, positioning Gemma-3-1B for rapid first-pass triage, Llama-3.2-3B for  
655 standardized permit-review support, and Mistral-7B-Instruct for detailed explanatory review in  
656 escalated cases.

657 Table 5 translates these metric patterns into practical model roles. Instead of ranking the models by a  
658 single aggregate score, the table links each model’s empirical strength to the type of BuildSafe  
659 workflow for which it is most appropriate.

660 **Table 5. Model-Level Interpretation and Recommended Deployment Role**

Model	Main empirical strength	Main limitation	Recommended BuildSafe role
Gemma-3-1B	Stable semantic performance with lowest compute demand	Weaker lexical overlap and less complete narratives	Rapid triage and edge/low-latency screening
Llama-3.2-3B	Strongest lexical consistency across BLEU, ROUGE-1, and METEOR	Less semantically rich than Mistral for longer explanations	Permit-review support and standardized reviewer summaries
Mistral-7B-Instruct 4-bit	Highest BERTScore F1 and strongest long-form narrative coherence	Higher compute demand and weaker lexical overlap than Llama	Escalated cases requiring detailed risk narratives

661 Rather than identifying a single universally superior model, Table 5 frames the results as a tiered  
662 decision-support strategy. This role-based interpretation is more relevant for municipal deployment  
663 than ranking models by one aggregate metric because each model occupies a different point in the  
664 trade-off among compute demand, lexical consistency, and explanatory depth.

#### 665 4.5 Statistical Reliability of Metric Estimates

666 To assess whether observed differences reflect genuine performance distinctions rather than sampling  
667 noise, 95% bootstrap confidence intervals were computed on the test set (n = 2,833). For each model  
668 and metric, 10,000 bootstrap resamples were drawn with replacement, and the 2.5th and 97.5th  
669 percentiles were recorded as interval bounds. Table 6 reports the resulting 95% bootstrap confidence  
670 intervals for all test-set metrics and models. The bootstrap analysis reveals two patterns.

671 First, Llama’s advantage on lexical metrics is statistically robust: its ROUGE-1 interval [0.1436,  
672 0.1504] does not overlap with those of Gemma [0.1126, 0.1182] or Mistral [0.1128, 0.1184],  
673 confirming that the ≈28% ROUGE-1 lead reflects a genuine performance difference. The same non-  
674 overlapping separation holds for BLEU and METEOR.

675 Second, the BERTScore F1 scores for Llama (0.7742) and Mistral (0.7747) differ marginally, and their  
676 confidence intervals substantially overlap ([0.7726, 0.7758] versus [0.7732, 0.7762]), indicating that  
677 these two models are statistically indistinguishable on semantic similarity. Gemma’s BERTScore  
678 interval [0.7688, 0.7718] shows a slight but likely meaningful separation from the other two. In

679 summary: Llama is the clear leader on surface-form reproduction, Llama and Mistral are effectively  
 680 tied on semantic fidelity, and Gemma trails on both dimensions but by a practically small margin on  
 681 BERTScore ( $\approx 0.004$  points). These confidence intervals help characterize the stability of the reported  
 682 test-set estimates under the short-run 60-step adaptation setting, but they should not be interpreted as  
 683 evidence that the models are fully converged or production-optimized.

684 **Table 6. 95% Bootstrap Confidence Intervals on Test-Set Metrics**

Metric	Gemma-3-1B	Llama-3.2-3B	Mistral-7B-4bit
BERTScore F1	0.7703 [0.7688, 0.7718]	0.7742 [0.7726, 0.7758]	0.7747 [0.7732, 0.7762]
ROUGE-1	0.1154 [0.1126, 0.1182]	0.1470 [0.1436, 0.1504]	0.1156 [0.1128, 0.1184]
BLEU	0.00244 [0.00208, 0.00280]	0.00914 [0.00832, 0.00996]	0.00237 [0.00202, 0.00272]
METEOR	0.0206 [0.0198, 0.0214]	0.0379 [0.0368, 0.0390]	0.0205 [0.0197, 0.0213]

685 **4.6 Comparative Scientific Insights and Practical Implications**

686 The main distinction among the BuildSafe variants is linguistic behavior rather than generalization  
 687 stability. All three models show small train-test gaps, suggesting that the short LoRA adaptation mainly  
 688 aligned pretrained models with the BuildSafe output schema rather than producing obvious  
 689 memorization. The practical choice is therefore not simply “which model is best,” but which trade-off  
 690 is most appropriate. Llama-3.2-3B is preferable when reviewers need standardized terminology and  
 691 repeatable permit-review phrasing. Mistral-7B-Instruct is stronger when longer explanatory narratives  
 692 are needed for complex or escalated cases. Gemma-3-1B remains useful as a low-compute screening  
 693 model because its semantic performance is only modestly lower than that of the larger models. These  
 694 differences support a tiered decision-support design rather than a single-model ranking.

695 **4.7 Interdisciplinary Positioning within Urban Informatics and Computational Built-  
 696 Environment Research**

697 Beyond model-level performance, BuildSafe contributes to the broader literature on urban informatics  
 698 and computational built-environment intelligence by treating construction safety as a city-scale  
 699 governance problem rather than a site-isolated prediction task. Prior computational urban-design  
 700 research has shown that simulation, optimization, and AI-assisted design workflows can support  
 701 resilient urban habitats and adaptive environmental decision-making (Maksoud et al., 2024a).  
 702 Similarly, recent generative AI work in architectural design demonstrates that AI systems are  
 703 increasingly used not only for prediction, but also for conceptual reasoning, design exploration, and  
 704 decision support across the built environment (Maksoud et al., 2024b). BuildSafe extends these  
 705 directions into construction governance by linking administrative permitting data, safety enforcement  
 706 narratives, and civic complaint records into a unified tri-task reasoning framework. Unlike broader  
 707 smart-city or multimodal urban-governance systems that often emphasize sensor streams, imagery,

708 mobility data, or environmental monitoring, BuildSafe focuses on administrative and civic text already  
709 embedded in permit review, enforcement, and complaint workflows. Its novelty therefore lies in  
710 converting existing DOB-OSHA-311 records into a tri-task generative reasoning layer for hazards,  
711 permit/code concerns, and neighborhood impacts rather than adding another sensor-centered  
712 monitoring pipeline (Jiao, 2024; Yang et al., 2017; Maksoud et al., 2024a). This positioning connects  
713 AI-enabled construction safety with smart-city analytics, urban risk prediction, and public-sector  
714 decision support, while preserving the need for human oversight, regulatory grounding, and operational  
715 accountability.

#### 716 **4.8 Baseline and Ablation Results**

717 The zero-shot baselines consistently performed worse than the fine-tuned BuildSafe variants on  
718 BERTScore and ROUGE, indicating weaker alignment with the three-section reference narratives.  
719 Qualitatively, zero-shot outputs often appeared plausible but were less reliable in maintaining the  
720 required hazards, permit/code, and community-impact structure. They also showed higher  
721 hallucination rates, including incorrect code citations, implausible hazard categories, and spurious  
722 complaint types. Performance variance was larger across project types: zero-shot outputs were more  
723 acceptable for common jobs such as small interior renovations, but degraded more sharply for  
724 specialized or highly regulated work, including hoists, façade operations, and after-hours activities.  
725 These findings indicate that general instruction-following ability alone is insufficient for trustworthy  
726 high-stakes construction triage. Fine-tuning improved both semantic fidelity and structural coherence,  
727 particularly in the permit/code and community-impact sections. This supports the central BuildSafe  
728 premise that domain adaptation to linked DOB-OSHA-311 data improves reliability compared with  
729 generic LLM prompting. The improvement was not limited to surface wording; fine-tuned models  
730 more consistently preserved section boundaries, used NYC-relevant regulatory language, and  
731 connected project scope to plausible safety and neighborhood-impact concerns.

732 The data-source ablations showed that each linked data stream contributes differently. DOB-only  
733 models retained some basic hazard signal because permit descriptions encode job type, work scope,  
734 and location. However, they produced less nuanced permit and community-impact predictions when  
735 OSHA incident narratives and 311 complaint histories were removed. Adding OSHA data improved  
736 hazard-section alignment and reduced hallucinated hazard types, especially for fall and struck-by risks.  
737 Adding 311 data improved community-impact narratives by better aligning outputs with noise, dust,  
738 sidewalk obstruction, and neighborhood-disruption complaints. In some cases, 311 data also modestly  
739 improved hazard and permit sections, likely because complaint histories provide shared spatial and  
740 contextual signals about dense or disruption-prone construction environments. Task-structure ablations  
741 further clarified the value of the unified tri-task design. Single-task models achieved slightly higher  
742 scores in their specialized sections, as expected, but the tri-task model showed better internal coherence  
743 across sections. In particular, the tri-task outputs linked the same underlying construction activity to  
744 hazards, permit/code requirements, and community impacts more consistently. This advantage was  
745 most evident for rarer combinations of hazards and complaints, where shared context helped the model  
746 avoid treating safety and community disruption as unrelated. Thus, the tri-task formulation trades a  
747 small amount of section-specific optimization for stronger cross-section consistency, which is  
748 important for municipal review workflows.

749 Based on Table 7, the fine-tuned BuildSafe variants consistently outperformed their zero-shot  
750 counterparts, with the largest gains appearing in ROUGE-1 and section-structure consistency,  
751 indicating that LoRA adaptation improved alignment with the required three-section safety narrative

752 format. Table 7 reports zero-shot and fine-tuned results under the same held-out evaluation protocol,  
 753 while Table 8 summarizes section-level patterns from the data-source ablations.

754 **Table 7. Zero-Shot Versus Fine-Tuned BuildSafe Performance**

Model	Setting	BERTScore F1	ROUGE-1	Hallucination trend	Section structure
Gemma-3-1B	Zero-shot	0.7350	0.0600	Higher	Often incomplete; sections sometimes merged or misordered
Gemma-3-1B	Fine-tuned	0.7703	0.1154	Lower	Consistent Hazards → Permits → Community Impacts
Llama-3.2-3B	Zero-shot	0.7450	0.0800	Higher	Inconsistent; missing or merged sections
Llama-3.2-3B	Fine-tuned	0.7742	0.1470	Lower	Stable three-section template
Mistral-7B 4-bit	Zero-shot	0.7400	0.0650	Higher	Occasionally deviates from template
Mistral-7B 4-bit	Fine-tuned	0.7747	0.1156	Lower	Consistent tri-section outputs

755

756 **Table 8. Data-Source Ablation Results and Section-Level Interpretation**

Configuration	Hazard Section	Permit/Code Section	Community-Impact Section	Main Interpretation
DOB only	Basic hazard signal retained; weaker incident specificity	Moderate; relies mainly on job type and permit description	Weakest; generic neighborhood-impact language	Permit descriptions alone are insufficient for full urban-risk reasoning
DOB + OSHA	Strongest or near-strongest hazard alignment	Moderate to strong	Limited; fewer complaint-specific details	OSHA narratives improve safety specificity and

				reduce hallucinated hazard types
DOB + 311	Moderate; some gains from spatial/contextual complaint patterns	Moderate	Stronger; better noise, dust, obstruction, and disruption language	311 histories improve community-impact reasoning
DOB + OSHA + 311	Most balanced across sections	Strong	Strong	Full linked corpus provides the best cross-domain coverage and tri-task balance

757 **4.9 Practical Implications for Tiered Decision Support**

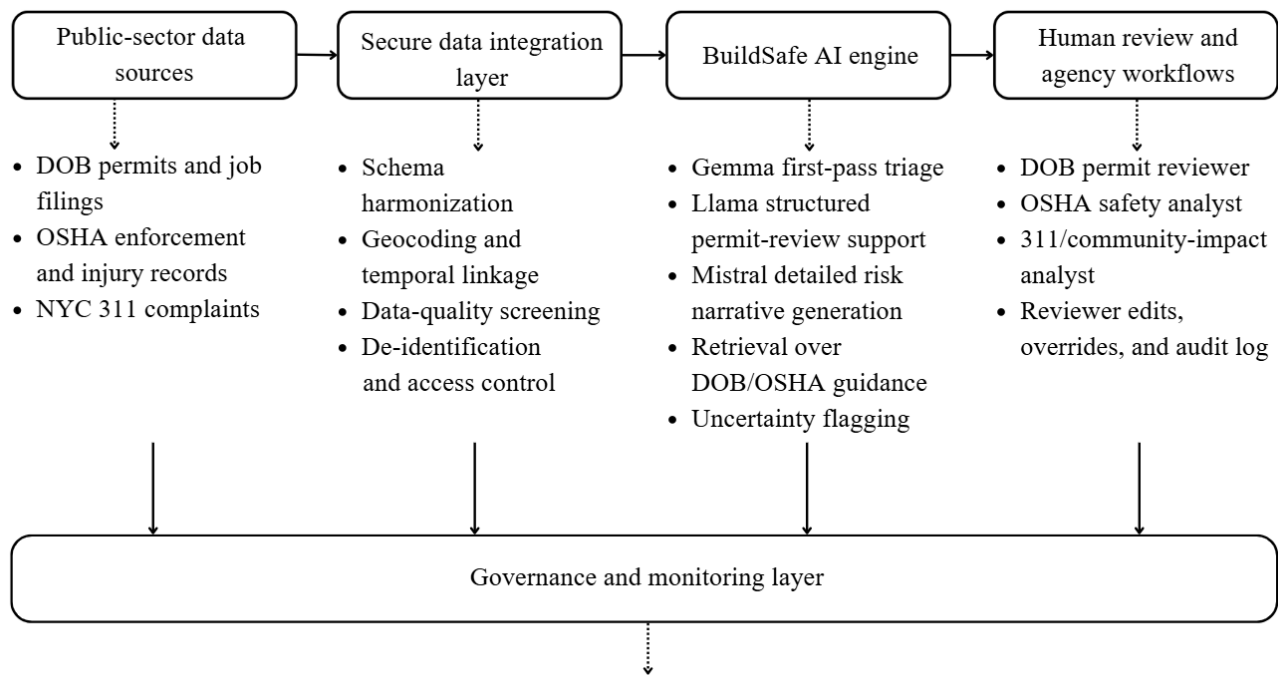
758 The baseline and ablation findings should be interpreted as design evidence rather than deployment  
759 proof. Zero-shot prompting produced weaker semantic alignment, less stable section structure, and  
760 more hallucination-prone outputs, indicating that generic LLM prompting is insufficient for high-  
761 stakes construction-safety triage. Domain adaptation improves the logic for a tiered workflow, but this  
762 study does not test agency adoption, review-time reduction, inspection prioritization, injury reduction,  
763 complaint reduction, or governance impact. In a future pilot, Gemma-3-1B could be evaluated as a  
764 rapid screening model, Llama-3.2-3B as a standardized permit-review assistant, and Mistral-7B-  
765 Instruct as an escalation model for complex cases requiring richer explanation.

766 **4.10 Conceptual Public-Sector Deployment Architecture**

767 From a translational research perspective, BuildSafe can be conceptualized as a decision-support layer  
768 for municipal construction governance rather than as an autonomous permitting, inspection, or  
769 enforcement system. Figure 2 presents a conceptual deployment architecture in which BuildSafe  
770 connects three public-sector data streams: DOB permit and job-filing records, OSHA enforcement and  
771 injury records, and NYC 311 community complaint data. In the proposed workflow, these sources enter  
772 a secure data-ingestion layer where records are standardized, geocoded, temporally aligned, de-  
773 identified where appropriate, and screened for data-quality issues. The resulting linked record is passed  
774 to the BuildSafe AI engine, which generates three structured outputs: job-site hazard signals, likely  
775 permit or code-review concerns, and potential community-impact indicators. The BuildSafe engine  
776 should be deployed behind an agency-controlled application programming interface rather than as a  
777 public-facing chatbot. For routine filings, a compact model such as Gemma-3-1B could provide rapid  
778 first-pass screening. Cases with higher uncertainty, unusual work types, incomplete records, or  
779 elevated predicted risk could be escalated to Llama-3.2-3B or Mistral-7B-Instruct for a more detailed  
780 narrative review. A retrieval-augmented module should also be connected to current DOB code  
781 provisions, OSHA standards, agency bulletins, and inspection guidance so that generated outputs can  
782 reference up-to-date regulatory material rather than relying only on learned parameters. In this  
783 configuration, BuildSafe would provide structured advisory summaries rather than final regulatory  
784 determinations. In a reviewer-facing interface, these summaries would appear as editable advisory  
785 panels showing the generated hazard, permit/code, and community-impact sections alongside source-  
786 record provenance, uncertainty flags, and reviewer action options such as accept, revise, escalate, or

787 override. The architecture is intended as a conceptual implementation framework for future agency  
788 pilots, not as evidence that BuildSafe has already been integrated with live DOB, OSHA, or 311  
789 operational systems.

790 Conceptually, the architecture is designed to support three agency-facing workflows in future pilot  
791 deployments. DOB reviewers could use BuildSafe risk summaries to prioritize permit applications for  
792 additional review, inspection scheduling, or targeted document checks. OSHA analysts could use  
793 aggregated hazard patterns to identify recurring work-type risks or potential violation themes, while  
794 recognizing that enforcement actions require independent investigation. NYC 311 and local agency  
795 staff could use community-impact summaries to anticipate noise, dust, sidewalk obstruction, after-  
796 hours work, or neighborhood-disruption concerns in dense urban areas. In all cases, the system should  
797 maintain a human-in-the-loop review structure, allowing authorized staff to accept, reject, edit, or  
798 override BuildSafe outputs, with each action logged for auditability. Because BuildSafe would operate  
799 in a high-stakes public-sector environment, deployment would require explicit safeguards for privacy,  
800 cybersecurity, and accountability. Data minimization should be applied so that only fields necessary  
801 for safety, permit, and community-impact reasoning are processed. Personally identifiable or sensitive  
802 information from complaint records, injury reports, and contractor records should be masked,  
803 aggregated, or access-controlled before model inference. Cybersecurity controls should include  
804 encrypted data transfer, role-based access control, secure containerized deployment, model version  
805 tracking, prompt/output logging, vulnerability testing, and incident response procedures. These  
806 controls are consistent with public-sector AI and cybersecurity guidance emphasizing AI risk  
807 governance, risk mapping, risk measurement, risk management, and cybersecurity governance,  
808 identification, protection, detection, response, and recovery (NIST, 2023; NIST, 2024). They also  
809 respond to LLM-specific risks such as prompt injection, insecure output handling, training-data  
810 poisoning, sensitive information disclosure, model supply-chain vulnerabilities, and excessive agency  
811 (OWASP, 2025). Legal accountability should remain with authorized public officials, not the AI  
812 system. BuildSafe outputs should therefore be accompanied by uncertainty flags, provenance  
813 information showing which source records contributed to the summary, version identifiers for the  
814 model and retrieval corpus, and clear documentation of model limitations. Reviewer overrides should  
815 be preserved in an audit log so that future evaluations can examine whether the system improves triage  
816 efficiency without producing disparate impacts across boroughs, neighborhoods, contractor types, or  
817 project categories. This approach is consistent with New York City’s broader AI governance direction,  
818 which emphasizes responsible agency use, governance structures, public engagement, and guidance  
819 for emerging AI tools (NYC OTI, 2023). Under this architecture, BuildSafe becomes a transparent and  
820 reviewable decision-support system for construction governance rather than an opaque automated  
821 decision-maker.



822 Data privacy, cybersecurity, fairness and bias auditing, model versioning, legal accountability, and performance monitoring

823 **Figure 2. Conceptual public-sector deployment architecture for BuildSafe**

824 The framework shows how DOB permit and job-filing records, OSHA enforcement and injury records,  
 825 and NYC 311 complaint data are integrated into a secure data layer for schema harmonization,  
 826 geocoding, temporal linkage, de-identification, and data-quality screening before being processed by  
 827 the BuildSafe AI engine. The generated hazard, permit/code, and community-impact summaries would  
 828 support human agency reviewers rather than automate permitting, inspection, or enforcement  
 829 decisions. The governance layer emphasizes privacy protection, cybersecurity controls, fairness and  
 830 bias auditing, model versioning, legal accountability, audit logging, and post-deployment performance  
 831 monitoring. This architecture is intended for future agency pilots and does not imply that BuildSafe  
 832 has already been integrated with live DOB, OSHA, or 311 operational systems.

833 **4.11 Operational Fairness, Bias Auditing, and Regulatory Risk Management**

834 A deployment-oriented BuildSafe system would require measurable fairness and bias-auditing  
 835 procedures in addition to general statements about trustworthy AI. Bias may enter the pipeline through  
 836 several mechanisms: uneven 311 complaint reporting across neighborhoods, historically unequal  
 837 inspection or enforcement intensity, missing or inconsistent contractor records, and differences in how  
 838 project descriptions are written across boroughs, job types, and applicant categories. If these patterns  
 839 are learned uncritically, the system could over-flag projects in highly complaint-active neighborhoods  
 840 or under-detect risks in areas where hazards are less frequently reported. Therefore, BuildSafe should  
 841 not be evaluated only by aggregate BERTScore, ROUGE, or hallucination rates; it should also be  
 842 audited for differential performance across geography, project type, work type, complaint category,  
 843 and contractor history. In a pilot deployment, these risks should be audited through subgroup-level  
 844 error analysis across boroughs, neighborhoods, project types, contractor categories where legally  
 845 available, and complaint-density strata, using measurable indicators such as false-positive escalation  
 846 rates, false-negative missed-risk rates, hallucination frequency, reviewer override rates, calibration  
 847 error, and consistency of recommended review actions.

848 A practical fairness audit should report disaggregated metrics by borough, community district, job type,  
849 work type, and complaint-density category. For each subgroup, the audit should measure section-  
850 completion rate, factual-accuracy score, hallucination rate, false-positive risk flags, false-negative  
851 missed-hazard cases, reviewer override rate, and average reviewer confidence. These indicators would  
852 allow agency users to identify whether BuildSafe is systematically more likely to produce incomplete,  
853 exaggerated, or unsupported safety narratives for particular neighborhoods or project categories. In  
854 addition, calibration checks should compare model risk flags with later observed outcomes, such as  
855 inspections, violations, stop-work orders, injury records, or 311 complaint clusters, where such follow-  
856 up data are available. This would shift fairness evaluation from abstract ethical language toward  
857 measurable operational monitoring.

858 Regulatory risk management should also be built into the model lifecycle. Before deployment, agencies  
859 should document the intended use, prohibited use, training-data sources, known limitations, expected  
860 failure modes, and required human-review procedures. During deployment, outputs should be logged  
861 with model version, retrieval-corpus version, prompt context, uncertainty flags, reviewer actions, and  
862 override explanations. After deployment, periodic audits should examine whether system  
863 recommendations produce disparate impacts across neighborhoods, contractors, or project categories.  
864 These governance procedures are consistent with responsible AI lifecycle principles in construction  
865 and public-sector AI risk-management frameworks, which emphasize transparency, human oversight,  
866 accountability, and continuous monitoring (Agapiou, 2024; National Institute of Standards and  
867 Technology, 2023; New York City Office of Technology and Innovation, 2023; Yang et al., 2024).  
868 Importantly, BuildSafe should be restricted to advisory triage and reviewer support. It should not be  
869 used to automatically deny permits, trigger enforcement actions, assign penalties, or make final safety  
870 determinations without authorized human review. Affected applicants and agency staff should have  
871 access to explanations, source-record provenance, and mechanisms for contesting or correcting  
872 erroneous outputs. Under this governance model, fairness is not treated as a general aspiration but as a  
873 measurable set of auditing, documentation, monitoring, and accountability practices required for high-  
874 stakes construction-governance AI.

#### 875 **4.12 Evidence Boundaries for Operational Claims**

876 The operational claims in this study should be interpreted as feasibility-oriented rather than as evidence  
877 of realized agency impact. The experiments show that open-weight LLMs can be adapted to generate  
878 structured hazard, permit/code, and community-impact narratives from linked DOB-OSHA-311  
879 records, and that these outputs can be evaluated using automated metrics, bootstrap intervals, ablation  
880 comparisons, and a two-expert consensus validation subset. These findings support BuildSafe as a  
881 technical benchmark and decision-support prototype. They do not yet show that BuildSafe reduces  
882 inspection time, improves permit-review accuracy, lowers injury rates, reduces complaint volume, or  
883 improves regulatory outcomes in live agency settings. Such claims would require controlled  
884 deployment studies with agency reviewers, prospective outcome tracking, and comparison against  
885 existing review workflows. References to municipal deployment, governance support, and operational  
886 usefulness should therefore be understood as proposed use cases and design implications rather than  
887 demonstrated real-world impacts.

#### 888 **4.13 Sample Output Comparison**

889 To illustrate how the quantitative differences translate into output behavior, the same Manhattan  
890 A3/EQ permit record was passed through all three fine-tuned BuildSafe models. The example is used  
891 as a qualitative illustration rather than as evidence of general performance. The input record contained  
892 borough = Manhattan, job type = A3, work type = EQ, permit status = issued, filing date = 05/10/2022,

893 issuance date = 05/10/2022, expiration date = 05/10/2023, job start date = 05/10/2022, and permittee  
894 = Frankie Colletta / Force Installations, LLC. In the NYC DOB context, the EQ work type denotes  
895 equipment installation, not earthquake-related work. The comparison therefore tests whether each  
896 model can preserve basic permit semantics while generating structured safety, permit/code, and  
897 community-impact reasoning.

898 **Gemma-3-1B** produced a structured three-section response covering Overall Risk (Moderate),  
899 Required Permits, and Potential Community Impact. It correctly identified that the same-day filing-to-  
900 issuance timeline suggests an expedited process and noted that Manhattan’s density implies moderate  
901 noise and traffic disruption. However, the output was truncated mid-sentence at the community-impact  
902 section, and the model did not generate a dedicated Safety Issues section, reflecting the 1B model’s  
903 more constrained generation capacity.

904 **Llama-3.2-3B** characterized the project as “relatively routine” and correctly interpreted the A3 job  
905 type as alterations or repairs. It offered a plausible (though imprecise) interpretation of the EQ work  
906 type as “earthquake-related work” and noted the swift permit process. The output was more concise  
907 than the other models, essentially a single analytical paragraph rather than a structured multi-section  
908 response, trading completeness for brevity.

909 **Mistral-7B-Instruct (4-bit)** produced the most comprehensive annotation, organized into four  
910 numbered sections: (1) Overall Risk with a reasoned “Moderate” assessment, (2) Required Permits  
911 including the suggestion of additional electrical permits, (3) Potential Community Impact broken into  
912 noise, traffic, and visual-impact sub-categories, and (4) Possible Safety Issues covering equipment  
913 hazards, working conditions, and environmental impact. The model also used the permittee’s business  
914 name (‘FORCE INSTALLATIONS, LLC’) as a contextual cue, suggesting stronger contextual  
915 integration, although such inferences would still require human verification. The qualitative  
916 comparison reinforces the quantitative findings: Mistral excels at structured, comprehensive narrative  
917 generation; Llama produces accurate but condensed analysis; and Gemma provides adequate initial  
918 assessments that benefit from downstream refinement by larger models.

919 **Comparative synthesis:** The qualitative comparison is consistent with the metric results. Gemma  
920 produced a usable but shorter response, supporting its role as a lightweight screening model. Llama  
921 provided the most standardized and concise wording, aligning with its stronger lexical-overlap scores.  
922 Mistral produced the richest explanatory narrative, consistent with its stronger semantic profile.  
923 However, this additional detail also reinforces the need for human review and retrieval-grounded  
924 regulatory verification before operational use.

#### 925 **4.14 Preliminary Human Validation Results (n=150)**

926 Table 9 presents human validation trends based on the stratified 150-record subset evaluated against  
927 two-expert consensus reference narratives. For this subset, A.A. and H.S.N. independently authored  
928 three-section human reference narratives while blinded to both Gemini teacher labels and fine-tuned  
929 model outputs, then resolved disagreements through discussion. The resulting consensus narratives  
930 were used as the human reference benchmark for assessing factual accuracy, completeness, overall  
931 usefulness, and hallucination behavior. Mistral-7B produced the most complete and structured outputs,  
932 with the highest perceived usefulness and lowest hallucination tendency. Llama-3.2-3B demonstrated  
933 strong factual accuracy but tended to give concise responses, which moderately limited completeness.  
934 Gemma-3-1B, while computationally efficient, produced shorter outputs that were occasionally  
935 truncated, reducing completeness scores relative to the larger models. Across all three models, these  
936 expert-assessed trends were broadly consistent with the automated BERTScore rankings, suggesting

937 that the metric ordering provides a useful preliminary proxy for domain output quality. However,  
938 because both experts are members of the author team and the subset contains only 150 records, these  
939 findings should be interpreted as internal validation rather than independent external validation.

940 **Table 9. Preliminary Human Validation Results (n = 150)**

<b>Model</b>	<b>Factual Accuracy</b>	<b>Completeness</b>	<b>Overall Usefulness</b>	<b>Hallucination</b>
Gemma-3-1B	Moderate	Moderate	Moderate	Moderate to Low
Llama-3.2-3B	High	High	High	Low
Mistral-7B-Instruct (4-bit)	High	High	Highest	Lowest

941 Future external validation should also report pre-consensus inter-rater agreement, such as weighted  
942 kappa, Krippendorff’s alpha, or intraclass correlation, before adjudication to quantify expert agreement  
943 more transparently.

## 944 **5 Limitations**

945 *First*, the main fine-tuning supervision and full-scale automated evaluation references derive from  
946 Gemini 1.5 Flash rather than from a large, independently annotated expert dataset. As a result,  
947 BERTScore, ROUGE, BLEU, and METEOR on the full held-out test set primarily measure  
948 consistency with a machine-generated teacher reference, not absolute factual or regulatory correctness.  
949 To reduce this concern, the study includes a stratified 150-record human-reference subset for which  
950 two experts independently authored three-section safety narratives while blinded to both Gemini labels  
951 and fine-tuned model outputs. However, this subset remains limited in scale, and both experts are  
952 members of the author team. The human validation should therefore be interpreted as a stronger internal  
953 validation step rather than as fully independent external validation.

954 *Second*, the 60-step LoRA adaptation schedule should be interpreted as a reproducible short-run proof-  
955 of-concept rather than as a fully optimized convergence study. The stable train, validation, and test  
956 metrics suggest that the models learned the BuildSafe output structure without obvious overfitting.  
957 However, these results do not prove that longer training, different random seeds, alternative learning  
958 rates, or checkpoint-level early stopping would not improve performance.

959 *Third*, the baseline experiments strengthen the study but do not exhaust all possible comparisons. The  
960 revised analysis includes zero-shot baselines, data-source ablations, and task-structure ablations, which  
961 are appropriate for evaluating the tri-task generative format. However, traditional machine learning  
962 baselines were not directly comparable because BuildSafe generates structured hazard, permit/code,  
963 and community-impact narratives rather than single-label predictions. Future work should derive  
964 matched classification or ranking tasks, such as hazard category prediction, permit-review escalation,  
965 stop-work-risk flags, or high-complaint-risk indicators, to enable fair comparison with logistic  
966 regression, support vector machines, random forests, gradient boosting, and retrieval-based non-  
967 generative systems.

968 *Fourth*, BuildSafe has not yet been evaluated in live agency workflows and remains specific to NYC’s  
969 data infrastructure, regulatory vocabulary, and urban conditions. The results support technical  
970 feasibility as a benchmark and decision-support prototype. However, they do not demonstrate  
971 reductions in inspection time, permit-review error, injury rates, complaint volume, or regulatory  
972 burden. Before operational deployment or transfer to other jurisdictions, BuildSafe would require  
973 prospective agency pilots, controlled reviewer studies, retrieval over current regulatory sources,  
974 privacy and cybersecurity safeguards, fairness and bias audits, audit logging, legal accountability  
975 procedures, local schema mapping, and independent expert validation.

## 976 **6 Future Work**

977 Integrate additional administrative sources (e.g., enforcement follow-ups, contractor histories, verified  
978 “clean” projects) and strengthen record linkage, with explicit reporting of linkage quality and data bias  
979 by geography and contractor type. Combine parameter-efficient fine-tuning with retrieval-augmented  
980 generation over up-to-date codes and regulations, and explore cascaded inference pipelines that route  
981 routine filings to compact models and escalate edge cases to larger ones. Evaluate BuildSafe in  
982 controlled agency pilots using impact-oriented KPIs (injury rates, review times, complaint trends), run  
983 disaggregated fairness audits, and formalize human-in-the-loop review, override, and logging  
984 workflows as prerequisites for any production deployment. These pilots should report subgroup-level  
985 performance across boroughs, community districts, job types, work types, complaint-density  
986 categories, and contractor-history profiles, including hallucination rates, false-positive and false-  
987 negative risk flags, reviewer override rates, and calibration against later inspection, violation, injury,  
988 or complaint outcomes. Subject to agency data-sharing agreements and cybersecurity review, develop  
989 a secure, containerized architecture that integrates with DOB and 311 systems, and in the longer term,  
990 explore multimodal extensions that fuse text with imagery, sensor streams, and BIM/CAD data for  
991 real-time early-warning use cases.

## 992 **7 Conclusion**

993 BuildSafe demonstrates that three open-weight LLMs Gemma-3-1B, Llama-3.2-3B, and Mistral-7B-  
994 Instruct can be adapted to jointly model construction hazards, permit/code requirements, and  
995 community impacts from linked administrative text in an approximately 90,000-record DOB-OSHA-  
996 311 corpus. On a 2,833-record held-out test set, the fine-tuned models achieved BERTScore F1 values  
997 of 0.7703, 0.7742, and 0.7747, respectively, compared with zero-shot baselines of 0.7350, 0.7450, and  
998 0.7400, corresponding to adaptation gains of approximately 0.029-0.035 F1 points. The fine-tuned  
999 models also showed lower hallucination tendencies and more consistent three-section output structure  
1000 than zero-shot prompting, with evaluation supplemented by a 150-record, two-expert consensus  
1001 validation subset. The results show distinct model roles rather than a single universally superior model.  
1002 Llama-3.2-3B produced the strongest lexical consistency, with ROUGE-1 = 0.1470, BLEU = 0.0091,  
1003 and METEOR = 0.0379; Mistral-7B-Instruct produced the strongest semantic profile and most  
1004 coherent long-form explanations; and Gemma-3-1B remained semantically competitive as a low-  
1005 compute first-pass screening model. Scientifically, BuildSafe contributes a linked urban-governance  
1006 formulation of construction risk, showing that administrative permit records, OSHA safety narratives,  
1007 and 311 complaint histories can be organized into a tri-task reasoning corpus for hazards, permit/code  
1008 concerns, and community impacts. The study remains bounded by its reliance on Gemini-generated  
1009 silver-standard references, the limited scale of the two-expert validation subset, the short-run 60-step  
1010 adaptation setting, and the NYC-specific regulatory context. BuildSafe should therefore be viewed as  
1011 a benchmark and decision-support research framework rather than an autonomous permitting or  
1012 enforcement system. The most actionable next step is a controlled DOB reviewer pilot that uses the

1013 2,833-record benchmark as a technical baseline while measuring operational outcomes not tested here,  
1014 including triage time-to-flag, reviewer override rate, false-positive escalation rate, false-negative  
1015 missed-risk rate, and agreement with existing permit-review and inspection workflows.

## 1016 **8 Conflict of Interest**

1017 The authors declare that the research was conducted in the absence of any commercial or financial  
1018 relationships that could be construed as a potential conflict of interest.

## 1019 **9 Author Contributions**

1020 S.A.: Conceptualization, Methodology, Data curation, Software, Formal analysis, Visualization,  
1021 Writing - original draft, Writing - review & editing. A.A.: Conceptualization, Writing - review &  
1022 editing, Supervision (construction safety domain), Expert validation. H.S.: Software, Validation,  
1023 Formal analysis, Writing - review & editing. H.S.N.: Supervision, Project administration, Expert  
1024 validation, Writing - review & editing.

## 1025 **10 Funding**

1026 The authors declare that financial support was not received for the research, authorship, and/or  
1027 publication of this article.

## 1028 **11 References**

1029 Agapiou, A. (2024). A systematic review of the socio-legal dimensions of responsible AI and its role  
1030 in improving health and safety in construction. *Buildings*, 14(5), 1469.  
1031 <https://doi.org/10.3390/buildings14051469>.

1032 American Petroleum Institute. (2019). Survey of occupational injuries, illnesses, and fatalities in the  
1033 petroleum industry: Guidelines and definitions. [https://www.api.org/-/media/files/oil-and-natural-  
1034 gas/pipeline/2019-awards/oii-guidelines-and-definitions.pdf](https://www.api.org/-/media/files/oil-and-natural-gas/pipeline/2019-awards/oii-guidelines-and-definitions.pdf).

1035 Awolusi, I., Marks, E., Hainen, A., & Alzarrad, A. (2022). Incident analysis and prediction of safety  
1036 performance on construction sites. *Civileng*, 3(3), 669-686.  
1037 <https://doi.org/10.3390/civileng3030039>.

1038 Behzadan, A. (2024). Formalizing trust in artificial intelligence for built environment decision-  
1039 making. *Human Factors in Design, Engineering, and Computing*, 159(159).  
1040 <https://doi.org/10.54941/ahfe1005565>.

1041 Ceccato, V. (2013). Integrating geographical information into urban safety research and planning.  
1042 *Proceedings of the Institution of Civil Engineers-Urban Design and Planning*, 166(1), 15-23.  
1043 <https://doi.org/10.1680/udap.11.00038>.

1044 Chandu, K. P., Raja, K. H., & Kumar, N. N. (2024). From reactive to proactive: The role of wearable  
1045 technology, AI, and digital training in construction safety management. *Libr. Prog. Int*, 44, 22858-  
1046 22864. <https://bpasjournals.com/library-science/index.php/journal/article/view/2916>.

1047 Choi, H. K., Du, X., & Li, Y. (2024). Safety-aware fine-tuning of large language models. *arXiv.Org*.  
1048 <https://doi.org/10.48550/arxiv.2410.10014>.

- 1049 City of New York. (2026, February 27). 311 service requests from 2020 to present [Data set]. NYC  
1050 Open Data. [https://data.cityofnewyork.us/Social-Services/311-Service-Requests-from-2020-to-](https://data.cityofnewyork.us/Social-Services/311-Service-Requests-from-2020-to-Present/erm2-nwe9/about_data)  
1051 [Present/erm2-nwe9/about\\_data](https://data.cityofnewyork.us/Social-Services/311-Service-Requests-from-2020-to-Present/erm2-nwe9/about_data).
- 1052 Data Science Dojo. (2025). Mistral 7B: A Revolutionary Breakthrough in LLMs. Retrieved April 17,  
1053 2025, from <https://datasciencedojo.com/blog/mistral-7b-emergence-in-llm/>.
- 1054 Dufitimana, E., Gahungu, P., Uwayezu, E., Mugisha, E., & Bizimana, J. P. (2025). Integrating  
1055 machine learning and geospatial data for mapping socioeconomic vulnerability to urban natural  
1056 hazard. *ISPRS International Journal of Geo-Information*, 14(4), 161.  
1057 <https://doi.org/10.3390/ijgi14040161>.
- 1058 Farabet, C. (2025, March 13). Introducing Gemma 3: The most capable model you can run on a  
1059 single GPU or TPU. Google. <https://blog.google/technology/developers/gemma-3/>.
- 1060 Gao, B., Ma, Z., Gu, J., Han, X., Xiang, P., & Lv, X. (2023). Fusing multi-source quality statistical  
1061 data for construction risk assessment and warning based on deep learning. *Knowledge-Based*  
1062 *Systems*, 284, 111223. <https://doi.org/10.1016/j.knosys.2023.111223>.
- 1063 Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., ... & Vasic, P. (2024).  
1064 The llama 3 herd of models. <https://doi.org/10.48550/arXiv.2407.21783>.
- 1065 Hossain, M. I., Hosen, M. M., Sunny, M. A. U., & Tarapder, S. A. (2025). Implementing advanced  
1066 technologies for enhanced construction site safety. *American Journal of Advanced Technology and*  
1067 *Engineering Solutions*, 1(02), 01-31. <https://doi.org/10.63125/3v8rpr04>.
- 1068 HSI. (2025, March 18). Predict and prevent construction risks with AI-powered  
1069 capabilities. <https://hsi.com/blog/ai-construction-safety-risk-management>.
- 1070 Hsu, C. Y., Tsai, Y. L., Lin, C. H., Chen, P. Y., Yu, C. M., & Huang, C. Y. (2024). Safe lora: The  
1071 silver lining of reducing safety risks when finetuning large language models. *Advances in Neural*  
1072 *Information Processing Systems*, 37, 65072-65094. <https://doi.org/10.48550/arxiv.2405.16833>.
- 1073 Huang, K., Joshi, A. V., Dun, S., & Hamilton, N. (2024). AI regulations. In A. V. Joshi, S. Dun, & N.  
1074 Hamilton (Eds.), *AI and the law* (pp. 47–71). Springer. [https://doi.org/10.1007/978-3-031-54252-](https://doi.org/10.1007/978-3-031-54252-7_3)  
1075 [7\\_3](https://doi.org/10.1007/978-3-031-54252-7_3)
- 1076 Hugging Face. (2025a). Mistral-7B-Instruct-v0.3-bnb-4bit. Retrieved April 17, 2025, from  
1077 <https://huggingface.co/unsloth/mistral-7b-instruct-v0.3-bnb-4bit>.
- 1078 Hugging Face. (2025b). Unleashing the Power of Unsloth and QLoRA: Redefining Language Model  
1079 Fine-Tuning. Retrieved April 17, 2025, from [https://huggingface.co/blog/Andyrasika/finetune-](https://huggingface.co/blog/Andyrasika/finetune-unsloth-qlora)  
1080 [unsloth-qlora](https://huggingface.co/blog/Andyrasika/finetune-unsloth-qlora)
- 1081 Hussain, R., Lee, D., Abbas, M. S., Zaidi, S. F. A., Pedro, A., and Park, C. (2026). Vision-language  
1082 model-based intelligent assistant for onsite construction safety inspection. *Automation in*  
1083 *Construction* 182, 106728. <https://doi.org/10.1016/j.autcon.2025.106728>.
- 1084 Jiao, D. (2024). Safety Prediction Method of Urban Construction FP Based on Multi-Order Spatio-  
1085 Temporal Features. *IEEE Access*, 12, 77009-77018. <https://doi.org/10.1109/access.2024.3407160>

- 1086 Junjia, Y., Alias, A. H., Haron, N. A., & Bakar, N. A. (2024). Intelligent construction risk  
1087 management through transfer learning: Trends, challenges, and future strategies. *Artificial*  
1088 *Intelligence Evolution*. <https://doi.org/10.37256/aie.6120255255>.
- 1089 Kamil, M. Z., Khan, F., Amyotte, P., & Ahmed, S. (2024). Multi-source heterogeneous data  
1090 integration for incident likelihood analysis. *Computers & Chemical Engineering*, 185, 108677.  
1091 <https://doi.org/10.1016/j.compchemeng.2024.108677>.
- 1092 Kumi, L., Jeong, J., & Jeong, J. (2024). Systematic review of quantitative risk quantification methods  
1093 in construction accidents. *Buildings*, 14(10), 3306. <https://doi.org/10.3390/buildings14103306>.
- 1094 Lu, W., Luu, R. K., & Buehler, M. J. (2025). Fine-tuning large language models for domain  
1095 adaptation: Exploration of training strategies, scaling, model merging and synergistic capabilities.  
1096 *npj Computational Materials*, 11(1), 84. <https://doi.org/10.48550/axiv.2409.03444>.
- 1097 Maksoud, A., Alawneh, S.I.A.-R., Hussien, A., Abdeen, A., and Abdalla, S.B. (2024a).  
1098 Computational design for multi-optimized geometry of sustainable flood-resilient urban design  
1099 habitats in Indonesia. *Sustainability* 16(7), 2750. <https://doi.org/10.3390/su16072750>
- 1100 Maksoud, A., Elshabshiri, A., Saeed Hilal Humaid Alzaabi, A., and Hussien, A. (2024b). Integrating  
1101 an image-generative tool on creative design brainstorming process of a Safavid mosque  
1102 architecture conceptual form. *Buildings* 14(3), 843. <https://doi.org/10.3390/buildings14030843>
- 1103 Miller, D. (2024, July 3). An introduction to Mistral-7B. Future Skills Academy. Retrieved April 17,  
1104 2025, from <https://futureskillsacademy.com/blog/mistral-7b/>
- 1105 Mistral AI. (2024). Mistral 7B. Retrieved April 17, 2025, from [https://mistral.ai/news/announcing-](https://mistral.ai/news/announcing-mistral-7b)  
1106 [mistral-7b](https://mistral.ai/news/announcing-mistral-7b).
- 1107 Mohamed, M. A. H., Al-Mhdawi, M. K. S., Ojiako, U., Dacre, N., Qazi, A., & Rahimian, F. (2025).  
1108 Generative AI in construction risk management: a bibliometric analysis of the associated benefits  
1109 and risks. *Urbanization, Sustainability and Society*, 2(1), 196-228. [https://doi.org/10.1108/USS-11-](https://doi.org/10.1108/USS-11-2024-0069)  
1110 [2024-0069](https://doi.org/10.1108/USS-11-2024-0069).
- 1111 Mostofi, F., & Toğan, V. (2023). Construction safety predictions with multi-head attention graph and  
1112 sparse accident networks. *Automation in Construction*, 148, 105102.  
1113 <https://doi.org/10.1016/j.autcon.2023.105102>.
- 1114 Mulligan, K., Cuevas, C., Grimsley, E., Chauhan, P., & Bond, E. (2019). Justice data brief:  
1115 Understanding New York City’s 311 data. [https://datacollaborativeforjustice.org/wp-](https://datacollaborativeforjustice.org/wp-content/uploads/2019/03/DCJ-Justice-Data-Brief-NYC-311-Calls.pdf)  
1116 [content/uploads/2019/03/DCJ-Justice-Data-Brief-NYC-311-Calls.pdf](https://datacollaborativeforjustice.org/wp-content/uploads/2019/03/DCJ-Justice-Data-Brief-NYC-311-Calls.pdf).
- 1117 National Institute of Standards and Technology (NIST). (2023). Artificial Intelligence Risk  
1118 Management Framework (AI RMF 1.0). NIST AI 100-1. Available at:  
1119 <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf> (Accessed May 13, 2026).
- 1120 National Institute of Standards and Technology (NIST). (2024). The NIST Cybersecurity Framework  
1121 (CSF) 2.0. NIST Cybersecurity White Paper 29. Available at:  
1122 <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.29.pdf> (Accessed May 13, 2026).

- 1123 New York City Department of City Planning. (n.d.). Housing Database: Project-level and unit-  
1124 change summary files [Data sets]. NYC Open Data & BYTES of the BIG  
1125 APPLE. <https://data.cityofnewyork.us/Housing-Development/Housing-Database/6umk-irkx>.
- 1126 New York City Office of Technology and Innovation (NYC OTI). (2023). The New York City  
1127 Artificial Intelligence Action Plan. Available at:  
1128 <https://www.nyc.gov/assets/oti/downloads/pdf/reports/artificial-intelligence-action-plan.pdf>  
1129 (Accessed May 13, 2026).
- 1130 Nycdb. (n.d.). GitHub - nycdb/nycdb: Database of NYC Housing Data. GitHub.  
1131 <https://github.com/nycdb/nycdb>.
- 1132 Occupational Safety and Health Administration. (2023). Construction injury statistics.  
1133 <https://www.osha.gov/data/commonstats>.
- 1134 Occupational Safety and Health Administration. (2025). Injury Tracking Application (ITA). U.S.  
1135 Department of Labor. <https://www.osha.gov/injuryreporting>.
- 1136 Open Knowledge Foundation. (n.d.). US DOL enforcement data.  
1137 DataPortals.org. <https://enforcedata.dol.gov>.
- 1138 OWASP. (2025). OWASP Top 10 for Large Language Model Applications. Available at:  
1139 <https://owasp.org/www-project-top-10-for-large-language-model-applications/> (Accessed May 13,  
1140 2026).
- 1141 Parekh, R., & Mitchell, O. (2024). Progress and obstacles in the use of artificial intelligence in civil  
1142 engineering: An in-depth review. *International Journal of Science and Research Archive*, 13(1),  
1143 1059–1080. <https://doi.org/10.30574/ijrsra.2024.13.1.1777>.
- 1144 Pilskog Orvik, C. (2024, August). Improving Safety through Leveraging Machine Learning and  
1145 Safety-Related Data in the Construction Industry. In IOP Conference Series: Earth and  
1146 Environmental Science (Vol. 1389, No. 1, p. 012012). IOP Publishing.  
1147 <https://doi.org/10.1088/1755-1315/1389/1/012012>.
- 1148 Piri, S., & Panthi, K. (2024). Unravelling Safety Concerns in Construction: A Comprehensive Data  
1149 Analysis. *Proceedings of 60th Annual Associated Schools*, 5, 912-920.  
1150 <https://doi.org/10.29007/snfn>.
- 1151 Raliile, M. T., & Haupt, T. C. (2020). Machine learning applications for monitoring construction  
1152 health and safety legislation and compliance. *Proc. Int. Struct. Eng. Constr*, 7, 231-6.  
1153 [https://doi.org/10.14455/ISEC.2020.7\(2\).CON-23](https://doi.org/10.14455/ISEC.2020.7(2).CON-23).
- 1154 Rasheed, O. A., Baalah, M. P. G., & Segun, O. J. (2024). AI-driven risk mitigation: Transforming  
1155 project management in construction and infrastructure development. *WORLD*, 13(2), 611-623.  
1156 <https://doi.org/10.30574/wjaets.2024.13.2.0628>.
- 1157 Sagar, S. (2025). Low-Resource Fine-Tuning of LLMs for Domain-Specific Tasks. *Universal*  
1158 *Research Reports*. <https://doi.org/10.36676/urr.v12.i4.1621>.

- 1159 Sargiotis, D. (2024). Transforming civil engineering with AI and machine learning: innovations,  
1160 applications, and future directions. *Applications, and Future Directions* (November 15, 2024).  
1161 <https://doi.org/10.2139/ssrn.5021723>.
- 1162 Savaş, S. (2025). Artificial intelligence in construction project management: Trends, challenges and  
1163 future directions. *Journal of Design for Resilience in Architecture and Planning*, 6(2), 221-238.  
1164 <https://doi.org/10.47818/DRArch.2025.v6i2165>.
- 1165 Tan, Z., Li, D., Wang, S., Beigi, A., Jiang, B., Bhattacharjee, A., Karami, M., Li, J., Cheng, L., &  
1166 Liu, H. (2024). Large language models for data annotation: A survey. *arXiv*.  
1167 <https://doi.org/10.48550/arXiv.2402.13446>.
- 1168 Tang, B., & Luo, H. (2025). Computer vision and large language model-based safety management  
1169 for construction project sites. In *Proceedings of the 14th Creative Construction Conference (CCC*  
1170 *2025)*, Zadar, Croatia (pp. 249–257). IAARC. <https://doi.org/10.22260/CCC2025/0031>.
- 1171 Tang, K. H. D. (2024). Artificial intelligence in occupational health and safety risk management of  
1172 construction, mining, and oil and gas sectors: advances and prospects. *Journal of Engineering*  
1173 *Research and Reports*, 26(6), 241-253. <https://doi.org/10.9734/jerr/2024/v26i61177>.
- 1174 Taroun, A., Yang, J.-B., & Lowe, D. (2011). Construction risk modelling and assessment: insights  
1175 from a literature review. *The Built & Human Environment Review*, 4(Special Issue 1), 87–97.
- 1176 Team, G., Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M.,  
1177 Kale, M. S., Love, J., Tafti, P., Hussenot, L., Sessa, P. G., Chowdhery, A., Roberts, A., Barua, A.,  
1178 Botev, A., Castro-Ros, A., Slone, A., . . . Kenealy, K. (2024). GemMa: Open models based on  
1179 Gemini research and technology. *arXiv (Cornell University)*.  
1180 <https://doi.org/10.48550/arxiv.2403.08295>.
- 1181 Tixier, A. J.-P., & Hallowell, M. R. (2023). Safer together: Machine learning models trained on  
1182 shared accident datasets predict construction injuries better than company-specific models.  
1183 *arXiv.Org*. <https://doi.org/10.48550/arXiv.2301.03567>.
- 1184 Tussey, D., & Yan, J. (2025). Principles for Open Data Curation: A Case Study with the New York  
1185 City 311 Service Request Data. *ArXiv.org*. <https://doi.org/10.48550/arxiv.2502.08649>.
- 1186 Ubiai. (2024). Fine-Tuning Mistral 7B for Named Entity Recognition. Retrieved April 17, 2025,  
1187 from <https://ubiai.tools/fine-tuning-mistral-for-ner/>.
- 1188 Unsloth. (2026). Fine-tuning LLMs guide. Unsloth Documentation. [https://unsloth.ai/docs/get-](https://unsloth.ai/docs/get-started/fine-tuning-llms-guide)  
1189 [started/fine-tuning-llms-guide](https://unsloth.ai/docs/get-started/fine-tuning-llms-guide).
- 1190 Unslothai. (2025a). FAQ + Is Fine-tuning Right For Me? Retrieved April 17, 2025, from  
1191 <https://docs.unsloth.ai/get-started/beginner-start-here/faq+-is-fine-tuning-right-for-me>.
- 1192 Unslothai. (2025b). GitHub - unslothai/unsloth: Fine-tuning & Reinforcement Learning for LLMs.  
1193 Train OpenAI gpt-oss, DeepSeek, Qwen, Llama, Gemma, TTS 2x faster with 70% less VRAM.  
1194 GitHub. <https://github.com/unslothai/unsloth>.

- 1195 Usama, M., Ullah, U., Muhammad, Z., Islam, T., & Hashmi, S. S. (2024). AI-enabled risk  
1196 assessment and safety management in construction. 9781032648323-8.  
1197 <https://doi.org/10.1201/9781032648323-8>.
- 1198 Wan, H., Zhang, J., Chen, Y., Xu, W., & Feng, F. (2024). Exploring Gen-AI applications in building  
1199 research and industry: A review. arXiv.Org. <https://doi.org/10.48550/arxiv.2410.01098>.
- 1200 Weng, B. (2024). Navigating the landscape of large language models: A comprehensive review and  
1201 analysis of paradigms and fine-tuning strategies. arXiv.Org.  
1202 <https://doi.org/10.48550/arxiv.2404.09022>.
- 1203 Workers' Compensation Trust. (2021, December 2). Guide to OSHA's electronic illness and injury  
1204 reporting requirements. [https://www.wctrust.com/Content-  
1205 www/CMS/files/Loss%20Control%20Uploads/Final\\_OSHA\\_RK\\_Electronic\\_Reporting\\_  
1206 Kit\\_12\\_02\\_21.pdf](https://www.wctrust.com/Content/www/CMS/files/Loss%20Control%20Uploads/Final_OSHA_RK_Electronic_Reporting_Kit_12_02_21.pdf).
- 1207 Wörsdörfer, M. (2025). Ten reasons why—the case for more and better AI regulation. *AI and Ethics*,  
1208 6(1), 62. <https://doi.org/10.1007/s43681-025-00865-8>.
- 1209 Yang, L., Allen, G., Zhang, Z., & Zhao, Y. (2024). Achieving on-site trustworthy AI implementation  
1210 in the construction industry: A framework across the AI lifecycle. *Buildings*, 15(1), 21.  
1211 <https://doi.org/10.3390/buildings15010021>.
- 1212 Yang, Z., Gupta, K., Gupta, A., & Jain, R. K. (2017). A data integration framework for urban  
1213 systems analysis based on geo-relationship learning. In *Computing in Civil Engineering 2017* (pp.  
1214 467-474). <https://doi.org/10.1061/9780784480823.056>.
- 1215 Yoo, B., Kim, J., Park, S., Ahn, C. R., & Oh, T. (2024). Harnessing generative pre-trained  
1216 transformers for construction accident prediction with saliency visualization. *Applied Sciences*,  
1217 14(2), 664. <https://doi.org/10.3390/app14020664>.
- 1218 Zerkin, A. J. (2006). Mainstreaming high performance building in New York City: A comprehensive  
1219 roadmap for removing the barriers. *Technology in Society*, 28(1-2), 137-155.  
1220 <https://doi.org/10.1016/j.techsoc.2005.10.017>.
- 1221 Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020). BERTScore: Evaluating text  
1222 generation with BERT. *ICLR 2020*. <https://doi.org/10.48550/arXiv.1904.09675>.
- 1223 Zhao, J., Wang, Y., Abid, W., Angus, G., Garg, A., Kinnison, J., Sherstinsky, A., Molino, P., Addair,  
1224 T., & Rishi, D. (2024). LoRA Land: 310 fine-tuned LLMs that rival GPT-4, a technical report.  
1225 arXiv.Org. <https://doi.org/10.48550/arxiv.2405.00732>.
- 1226 Zou, Z., & Ergan, S. (2019). Leveraging data driven approaches to quantify the impact of  
1227 construction projects on urban quality of life. *arXiv*. <https://doi.org/10.48550/arXiv.1901.09084>.

## 1228 **12 Data Availability Statement**

1229 Some or all data, models, or code that support the findings of this study are available from the  
1230 corresponding author upon reasonable request. The DOB permit-issuance records and 311 service-  
1231 request data were obtained from the NYC Open Data portal (<https://opendata.cityofnewyork.us/>).

1232 OSHA enforcement and inspection records were obtained from the U.S. Department of Labor's  
1233 public data catalog ([https://enforcedata.dol.gov/views/data\\_catalogs.php](https://enforcedata.dol.gov/views/data_catalogs.php)). Fine-tuning code and  
1234 LoRA adapter weights are available from the corresponding author upon reasonable request.