

Causal-Constrained GraphRAG: Enhancing Explainability and Reliability in Safety-Critical Parameter Extraction

Anonymous ACL submission

Abstract

Causal-Constrained GraphRAG is a neurosymbolic framework designed to improve the explainability and reliability of parameter extraction from long, structured, and causally rich regulatory corpora used in safety-critical engineering domains such as fire safety and evacuation modeling. The framework parses heterogeneous technical manuals into a topologically grounded knowledge graph, extracts explicit cause-effect-action chains through weak supervision combined with an active-learning annotation loop, and injects validated causal structures into a GraphRAG reasoning stage to constrain large language model inference. Instantiation on a curated Evacuation and Fire Safety question-answer dataset (approximately 24k pairs) and a Neo4j-based causal knowledge graph enables evaluation across multi-hop regulatory reasoning, numerical and sign consistency, and end-to-end parameter generation tasks. Experimental results demonstrate substantial improvements in factual retrieval accuracy and causal coherence, including notable F1 gains in occupant-load interpretation and high mean opinion scores from domain experts, alongside improved response stability under repeated technical queries. The framework supports safe and auditable automation of engineering parameter extraction while explicitly preserving human-in-the-loop validation for cases involving ambiguous, incomplete, or conflicting regulatory guidance.

1 Introduction

The deployment of Large Language Models (LLMs) in safety-critical engineering domains, such as fire safety and evacuation modeling, is hindered by a persistent reliability gap. Standard Retrieval-Augmented Generation (RAG) frameworks have demonstrated proficiency in factual retrieval, but they frequently struggle with the implicit reasoning required to navigate complex regulatory landscapes (Gwynne et al., 2015). In these

domains, practitioners must synthesize information from fragmented technical manuals, such as NFPA 101 (Tubbs et al., 2024), BS 9999 (British Standards Institution, 2017), and the SFPE Handbook (Hurley et al., 2016), to derive precise and mutually consistent simulation parameters.

A critical bottleneck arises when simulation-ready parameters are missing or underspecified within the source text (Dimyadi et al., 2008). Traditional LLMs, when faced with such gaps, are prone to hallucinations or rely on shallow statistical correlations rather than the domain-specific causal logic required for safety-critical validation (Durmus et al., 2024). For instance, a model may successfully retrieve a baseline walking speed but fail to adjust it for the interdependent causal effects of smoke density and visibility unless explicitly guided by a structured knowledge framework (Kuligowski and Gwynne, 2010). To address these challenges, we propose a *Neurosymbolic framework* that integrates the linguistic flexibility of LLMs with the structural rigor of Causal-Constrained GraphRAG. Our approach moves beyond simple vector-based retrieval by parsing unstructured regulations into a topologically grounded knowledge graph (Edge et al., 2024). Furthermore, we introduce an NLP pipeline that leverages weak supervision and active learning to extract explicit causal chains, such as *SmokeDensity* $\uparrow \rightarrow$ *Visibility* $\downarrow \rightarrow$ *WalkingSpeed* \downarrow , directly from technical corpora.

By injecting these causal structures into the reasoning process of LLMs and a fine-tuned Gemma 3 (4B) model, we ensure that inferred parameters are both explainable and compliant with established engineering standards. Our evaluation, including multi-hop reasoning tasks and Mean Opinion Scores (MOS) from domain experts, demonstrates that this framework significantly improves the causal coherence and technical accuracy of pa-

parameter extraction in high-stakes environments.

Although we instantiate and evaluate our framework in fire safety and evacuation modeling, the underlying methods generalize to other regulatory domains involving long, structured, and causally rich corpora (e.g., transportation safety guidelines, accessibility standards, and healthcare protocols).

Contributions

- ✍ We formulate causally-constrained regulatory parameter inference as an NLP task and construct a regulatory QA resource over fire and evacuation standards to support training and evaluation.
- 🧠 We propose a Causal-Constrained GraphRAG architecture that combines graph-structured retrieval over regulatory codes with a neurosymbolic causal knowledge graph extracted via weak supervision and active learning.
- 🛡 We conduct a comprehensive evaluation of multiple LLMs (including a fine-tuned 4B-parameter model) on multi-hop regulatory QA, causal reasoning, and end-to-end parameter generation, with expert MOS and stability analyses in a safety-critical setting.

2 Related Work

2.1 RAG in Technical and Safety-Critical Domains

Traditional RAG systems primarily utilize vector-only retrieval, which excels at surface-level factual recall but often fails to capture the structural logic embedded in regulatory codes (Gao et al., 2024). In safety-critical fields like fire engineering, parameters are rarely standalone; they are governed by interdependent constraints found in manuals such as NFPA 101 and BS 9999. Previous research indicates that manual extraction from these sources is highly prone to error and inconsistency (Gwynne et al., 2015; Rafe et al., 2025). Emerging AI-driven approaches have attempted to automate this via NLP parsing and LLM-based retrieval (Li et al., 2024; Madireddy et al., 2025), but they typically treat parameters in isolation and do not jointly reason over interacting quantities such as occupant load, exit capacity, and smoke conditions.

Moreover, most existing RAG pipelines for technical domains implement a flat, chunk-based re-

trieval mechanism that ignores the underlying document structure and the logical relationships among provisions (Sarathi et al., 2024; Huang et al., 2025; Lu et al., 2025). As a result, they can retrieve relevant fragments while still producing parameter sets that violate implicit constraints in the standards (e.g., consistency between egress width, occupant density, and allowable evacuation times). In contrast, GraphRAG-style methods that operate over knowledge graphs promise more structure-aware retrieval (Edge et al., 2024), yet they have rarely been specialized to safety codes or tightly integrated with explicit causal constraints. Our work builds on this line by constructing a graph representation of regulatory sources and augmenting it with causal relations tailored to evacuation parameters.

2.2 Causal NLP and Knowledge Extraction

A major limitation of current LLMs is their reliance on statistical correlation rather than causal logic, which can result in counterintuitive or unsafe parameter sets in safety-critical settings (Hashmi et al., 2024). Recent advancements in Causal NLP have focused on extracting causal graphs from unstructured text and using them to guide model behavior (Wang et al., 2025; Kiciman et al., 2023). For example, Kiciman et al. (2023) demonstrate that causal-aware LLMs significantly outperform correlation-based systems in pairwise inference tasks when provided with explicit causal structures. Complementary work shows that embedding knowledge graphs into the reasoning process can improve generalization under limited or noisy data (Kiciman et al., 2023; Luo et al., 2024; Wu and Tsioutsoulouklis, 2024).

However, most existing causal-NLP frameworks are evaluated on short sentences, news-style text, or synthetic benchmarks, and they rarely address long, highly structured regulatory documents with dense cross-references. Furthermore, prior work typically focuses on causal inference tasks in isolation rather than on end-to-end generation of domain-specific artifacts such as simulation parameter sets.

3 Methodology

We propose a modular Neurosymbolic pipeline designed to transform unstructured regulatory corpora into causally-coherent simulation parameters. The framework consists of three primary stages: (1) structured multi-modal parsing, (2) targeted causal relationship extraction, and (3) Causal-Constrained

GraphRAG inference.

3.1 Multi-Modal Document Parsing and Structuring

Technical manuals such as NFPA 101 and the SFPE Handbook contain dense tabular data and hierarchical sectioning often lost in standard PDF-to-text conversion. To preserve semantic integrity, we evaluated multiple parsers using the DP-Bench benchmark (Upstage, 2024). LlamaParse (Liu, 2022) was selected as the primary engine, achieving a Normalized Indel Distance (NID) of 98.54 and a Tree Edit Distance-based Similarity (TEDS) of 92.38 for table recognition.

3.2 Targeted Causal Relationship Extraction

To address missing parameters, we implement a seven-stage extraction pipeline to identify "cause-effect-action" structure:

1. **Data Collection and Preprocessing:** Textual content from authoritative technical manuals, including the SFPE Handbook and NFPA 101—is cleaned and segmented into a sentence-level dataset $S = \{s_1, s_2, \dots, s_N\}$.

2. **Candidate Identification:** We utilize dependency parsing via spaCy (Honnibal et al., 2020) to identify marked causality (e.g., "results in") and unmarked causality (e.g., implied relations).

3. **Weak Supervision:** Using the Snorkel framework (Ratner et al., 2017), we apply labeling functions Λ to assign probabilistic labels. The confidence score is computed as:

$$P(\text{causal}|s_l, \Lambda) = \frac{\sum_{j=1}^m w_j P(\lambda_j(s_l)|\text{causal})}{\sum_{j=1}^m w_j} \quad (1)$$

where w_j is a weight assigned to trigger-specific labeling functions (e.g., $w_j = 1.5$ for simulation variables).

4. **Manual Labeling:** Top candidates are manually annotated for cause/effect spans and direction $d \in \{+, -\}$ using domain expertise[cite: 310, 312].

5. **Active Learning Loop:** We employ entropy-based uncertainty sampling via modAL (Danka and Horvath, 2018). Uncertainty $\mathcal{H}(s_i)$ is calculated using the multiline approach to fit column width:

$$\begin{aligned} \mathcal{H}(s_i) &= -P(\text{causal}|s_i) \log P(\text{causal}|s_i) & 208 \\ &- (1 - P(\text{causal}|s_i)) \log(1 - P(\text{causal}|s_i)) & 209 \\ & & (2) \end{aligned}$$

where the causal probability follows a Logistic Regression model[cite: 320]:

$$P(\text{causal}|s_i) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{in})}} \quad (3) \quad 212$$

6. **Structure Inference:** Extracted relations are categorized by a reasoning model into chains ($v_1 \rightarrow v_2 \rightarrow v_3$), forks ($v_1 \leftarrow v_2 \rightarrow v_3$), or colliders ($v_1 \rightarrow v_2 \leftarrow v_3$).

7. **Graph Database Storage:** Finally, the structured causal relationships are encoded into a Neo4j (Neo4j, Inc., 2025) graph database, where each node represents a simulation variable and directional edges express validated causal influences.

3.3 Causal-Constrained GraphRAG and Reasoning Pipeline

The final stage implements a GraphRAG architecture to bridge the gap between static regulatory knowledge and dynamic reasoning for parameter extraction. This system utilizes a Neo4j knowledge graph to provide a topologically grounded reference for the language models.

Graph Construction The Neo4j Knowledge Graph is constructed by encoding structured document segments and extracted causal structures into a property graph $G = (V, E)$. Nodes V represent simulation variables (e.g., *Smoke Density*, *Walking Speed*), while directed edges E represent validated causal influences (chains, forks, or colliders) or regulatory hierarchies. For instance, the graph for NFPA 101 encompasses over 33,000 nodes and 76,000 relationships.

Causal-Aware "Thinking" Prompts To fill parameter gaps where direct manual guidance is absent, we employ a "Thinking-Mode" reasoning framework. The process utilizes a structured four-step prompt to force the model to reason through the causal chain before delivering a final numerical estimate.

As defined in the prompt structure, the model is explicitly instructed to align its reasoning with retrieved causal relationships by injecting validated paths (e.g., *Smoke Density* $\uparrow \rightarrow$ *Visibility* $\downarrow \rightarrow$

252	<i>Speed</i> ↓) into the context, thereby constraining	
253	inference within physical laws and domain logic.	
254	Reasoning is grounded in structured baseline val-	
255	ues extracted from authoritative standards such as	
256	SFPE and BS 9999, ensuring regulatory consis-	
257	tency. When direct numerical guidance is unavail-	
258	able, the model analyzes underlying causal mech-	
259	anisms to derive reasonable estimates using logi-	
260	cal constraints and appropriate statistical assump-	
261	tions, such as log-normal distributions for delay	
262	processes. Each inferred output is accompanied	
263	by a structured justification to verify adherence to	
264	established domain dependencies and regulatory	
265	constraints.	
266	4 Experimental Setup	
267	4.1 Dataset and Corpus	
268	The experimental corpus is built upon authorita-	
269	tive fire safety and evacuation manuals, including	
270	NFPA 101, BS 9999, and the SFPE Handbook.	
271	From these sources, we curated the <i>Evacuation</i>	
272	<i>and Fire Safety Q&A Dataset</i> , comprising approxi-	
273	mately 24,000 question-answer pairs. The dataset	
274	was expanded using a three-stage process: manual	
275	extraction by domain experts (D_0), LLM-assisted	
276	expansion (D_{LLM}), and entropy-based active learn-	
277	ing refinement (D_{AL}) to optimize for edge-case	
278	coverage and technical density.	
279	4.2 Model Selection and Fine-Tuning	
280	We evaluated four primary model configurations,	
281	including commercial reasoning models, GPT-5.2,	
282	Claude Opus 4.5, and Gemini 3 Pro, as well as	
283	a fine-tuned domain-specific model based on	
284	Gemma 3 (4B) trained on the curated Q&A dataset	
285	described earlier. Fine-tuning of Gemma 3 was con-	
286	ducted using the Unsloth framework on a single	
287	NVIDIA A100-SXM4 GPU with mixed-precision	
288	(bfloat16) training. Optimization minimized cross-	
289	entropy loss over the task-adapted parameters θ^* ,	
290	converging to a stable minimum loss of 1.20 after	
291	approximately 400 training steps.	
292	4.3 Evaluation Tasks and Metrics	
293	Models were benchmarked across four distinct	
294	safety-critical task categories:	
295	1. <i>Tabular Data Interpretation</i> : Extracting spe-	
296	cific factors (e.g., occupant load factors) from	
297	structured regulatory tables.	
	2. <i>Textual Information Retrieval</i> : Accurate	298
	clause identification and retrieval from un-	299
	structured text.	300
	3. <i>Numerical Sign/Units Handling</i> : Correct de-	301
	tection and interpretation of signs (e.g., nega-	302
	tive tolerances) and unit conversions.	303
	4. <i>Calculation Application</i> : Multi-step applica-	304
	tion of regulatory formulas (e.g., exit capacity	305
	based on area and density).	306
	Performance was quantified using standard	307
	IR/QA metrics: Precision, Recall, and F1-score.	308
	For qualitative reasoning and causal consistency,	309
	we utilized the MOS) scale (1–5), where domain	310
	experts assessed technical accuracy, standards com-	311
	pliance, and causal coherence.	312
	4.4 Case Study: Spatial Layout Injection	313
	To evaluate the framework’s end-to-end perfor-	314
	mance in a real-world setting, we utilized a four-	315
	story educational building as a case study. This	316
	multi-functional facility, comprising classrooms,	317
	laboratories, administrative offices, and a public	318
	café, provided a complex architectural environment	319
	to test spatial grounding and functional context.	320
	We transformed building plans into a topologi-	321
	cal graph $G_{spatial} = (V, E)$. In this representation,	322
	nodes V denote individual rooms or circulation	323
	zones, while edges E correspond to horizontal or	324
	vertical connections such as corridors, doors, and	325
	stairwells. Each node was enriched with geometric	326
	and semantic attributes, including net area and oc-	327
	cupancy classifications. This spatial graph was in-	328
	tegrated into the RAG pipeline, allowing the LLMs	329
	to query and generate simulation-ready parameters	330
	(e.g., occupant load) tailored to the specific func-	331
	tional constraints of each individual room.	332
	5 Results and Analysis	333
	5.1 Retrieval Performance and Task Accuracy	334
	The retrieval accuracy was assessed across four	335
	technical categories as illustrated in Figure 1.	336
	Claude Opus 4.5 consistently achieved the highest	337
	performance, with F1 scores ranging from 0.917 to	338
	0.937. Specifically, it reached a peak F1-score of	339
	0.937 in <i>Textual Information Retrieval</i> , demonstrat-	340
	ing a superior capacity for long-context semantic	341
	parsing and clause resolution in regulatory corpora.	342

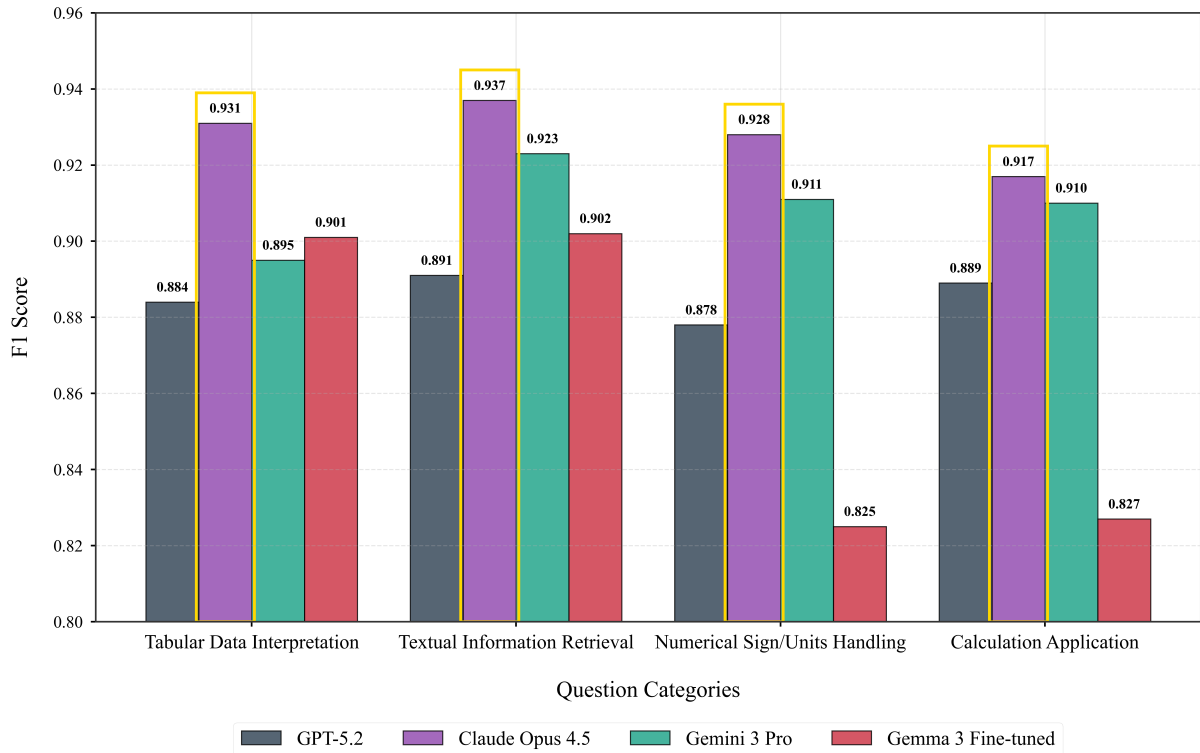


Figure 1: Performance of four LLMs across evacuation-related question categories.

The Specialization vs. Scale Trade-off The Gemma 3 Fine-tuned model (4B) exhibited significant adaptability, achieving an F1 score of 0.902 in *Textual Information Retrieval*. As illustrated in Figure 2, this reflects substantial domain-specific improvements across all regulatory categories compared to the base model.

The most pronounced advancement was observed in *Occupant Load* interpretation, which saw a +140.2% improvement in BERT-F1 score (increasing from 0.390 to 0.937). Similarly, *Fire Protection* and *Emergency Planning* tasks yielded gains of +77.2% and +16.0% respectively, validating the efficacy of the three-stage Q&A dataset generation process.

However, a performance ceiling was observed in computational tasks like *Numerical Sign/Units Handling* (0.825) and *Calculation Application* (0.827) as shown in Figure 1. From an NLP perspective, this suggests that while fine-tuning on domain-specific Q&A successfully encodes *lexical* and *terminological* knowledge, the *symbolic reasoning* required for mathematical manipulation remains highly dependent on the pre-training scale of larger models like Claude Opus 4.5. This validates the hybrid architecture of our framework: utilizing a fine-tuned small-language model (SLM)

for efficient textual classification while delegating complex multi-step arithmetic to higher-capacity reasoning engines.

5.2 Reasoning Capabilities and Stability

A critical metric for safety-critical NLP is the ability to maintain consistency across repeated technical queries, a prerequisite for engineering reproducibility. In our stability tests (Table 1), Claude Opus 4.5 demonstrated high reliability, returning stable responses in approximately 75% – 80% of trials.

Impact of Document Heterogeneity A significant finding emerged regarding the linguistic structure of source documents. Gemma 3 Fine-tuned achieved perfect stability (10/10) with the structured, rule-heavy content of *BS 9999*. However, its consistency plummeted to 5/10 when faced with *NFPA 101*, a document characterized by greater heterogeneity, recursive cross-references, and implicit clause logic. This indicates that SLMs, even when fine-tuned, are more susceptible to "contextual noise" and "distractor clauses" within complex technical documents. In contrast, the "thinking mode" in Claude Opus 4.5 acts as a self-correction mechanism, allowing it to navigate heterogeneous layouts with consistent attention to the relevant

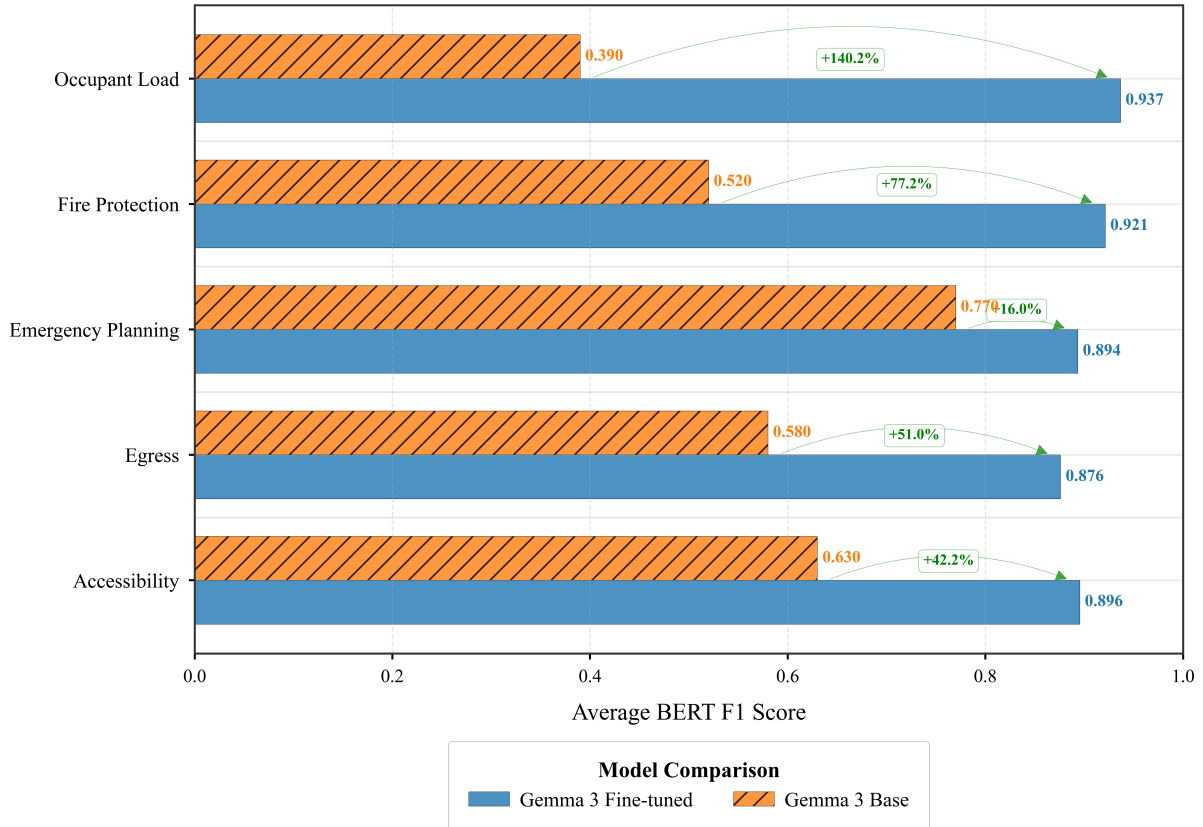


Figure 2: Comparative performance gains of Gemma 3 Fine-tuned over Gemma 3 Base across safety-critical domain categories.

LLM	SFPE Handbook	BS 9999	NFPA 101
Claude Opus 4.5	8/10	7/10	8/10
Gemma 3 (4B) Fine-tuned	7/10	10/10	5/10

Table 1: Stability Test Results: Number of stable, identical responses out of 10 queries per source.

topological links in the Neo4j graph.

5.3 Causal-Constrained Multi-hop Reasoning

Evaluation of multi-hop reasoning performance (Table 2) reveals the limitations of standard statistical retrieval. Claude Opus 4.5 successfully integrated complex, multi-variable causal chains (e.g., *smoke* → *visibility* → *speed*) in 9 out of 10 cases.

Neurosymbolic Injection vs. Stochastic Generation Most general-purpose models (GPT-5.2, Gemini 3 Pro) displayed a tendency to oversimplify or truncate causal logic when document segments were non-contiguous. This "shallow reasoning" is a byproduct of pure statistical prediction. By injecting the Causal Relationship Graph (Figure 3-10) directly into the prompt (Code 3-1), we force the model to adopt a *Neurosymbolic* reasoning path.

Instead of predicting the next likely token based on correlation, the model must traverse the chain $v_1 \rightarrow v_2 \rightarrow v_3$ (e.g., *Smoke* ↑ → *Visibility* ↓ → *Speed* ↓) before arriving at a numerical estimate. This ensures that "missing" parameters are not just guessed, but inferred via physical and regulatory logic.

5.4 Explainability and Expert Evaluation (MOS)

The core contribution of this work is "Mechanistic Interpretability"—the ability to show *why* a parameter was chosen. This was validated through expert MOS in Table 3.

Experts rated Claude Opus 4.5 with a near-perfect *Causal Coherence* score of 4.9. The qualitative feedback indicated that while models like GPT-5.2 were efficient, their lack of explicit causal

Model	Accuracy	Linguistic & Reasoning Performance Notes
Claude Opus 4.5	9/10	Successfully integrated multi-variable causal chains with interpretable justifications and realistic statistical distributions.
GPT-5.2	8/10	Provided high numerical precision but occasionally exhibited "hallucinated simplifications" where complex causal pathways were truncated.
Gemini 3 Pro	7/10	Strong mathematical derivations; however, the model often prioritized numeric modeling over behavioral logic.
DeepSeek R1	7/10	Precise outputs for specific variables, but exhibited "weak contextual transition" when reasoning across three or more interdependent steps.
Sonar Reasoning Pro	6/10	Heavily grounded in citations, but failed to synthesize multi-hop logic when parameters were spread across non-contiguous document sections.

Table 2: Multi-hop reasoning accuracy and performance notes for each LLM configuration.

Model	Tech. Accuracy	Std. Compliance	Causal Coherence	Avg. MOS
Claude Opus 4.5	4.8	4.5	4.9	4.7
GPT-5.2	4.5	4.5	4.4	4.5
Gemini 3 Pro	4.6	4.2	3.7	4.2
Sonar Reasoning Pro	4.2	4.4	3.8	4.1
DeepSeek R1	3.5	4.4	4.0	4.0

Table 3: Mean Opinion Score (MOS) for LLM outputs by domain experts.

constraints led to parameter sets that were occasionally technically "correct" but logically disjointed. Our framework's output, by contrast, was deemed "trustworthy for engineering use" because the justification explicitly mirrors the causal logic found in authoritative manuals like the *SFPE Handbook*.

5.5 Generalization to Safety-Critical Domains

The primary implication of these results is that parameter extraction in safety-critical domains cannot rely solely on larger datasets. The +140.2% improvement in the fine-tuned Gemma model proves that domain-specific alignment is essential, but the stability and multi-hop reasoning results prove that an external "causal guardrail" (GraphRAG) is necessary to ensure those models operate within the bounds of physical and regulatory laws. This framework provides a reproducible blueprint for other high-stakes domains, such as medical protocol extraction or legal compliance, where "shallow retrieval" is insufficient.

6 Conclusion

This study introduced a *Neurosymbolic framework* to address the fundamental reliability gap in stan-

dard RAG pipelines for safety-critical engineering domains. By integrating a Causal-Constrained GraphRAG architecture with a 7-stage causal extraction pipeline, we demonstrated a robust methodology for transforming unstructured regulatory documents into verifiable simulation parameters.

Our findings establish that injecting explicit causal structures, such as chains, forks, and colliders, significantly bounds the "stochastic generation" of LLMs within established physical and regulatory logic. The empirical evaluation showed that high-capacity reasoning models like Claude Opus 4.5 achieved superior accuracy in multi-hop reasoning tasks, while our domain-specialized Gemma 3 Fine-tuned model demonstrated a improvement in interpreting specialized occupant load provisions. The successful generation of a simulation-ready XML file for a complex educational building case study validates the framework's readiness for real-world engineering workflows.

The primary implication of this work is a major shift from subjective, manual parameter extraction toward an objective, transparent, and reproducible process. This approach not only enhances the safety of performance-based design but also

478	provides a generalizable blueprint for other high-		
479	stakes domains (e.g., medical and legal compli-		
480	ance) where hallucinatory outputs are unaccept-		
481	able.		
482	Limitations		
483	Despite its performance, the framework’s output		
484	quality remains inherently dependent on the clarity		
485	and comprehensiveness of the source regulatory		
486	documents. While the causal-constrained reason-		
487	ing mitigates hallucinations, the system is not in-		
488	fallible and still requires human-in-the-loop expert		
489	oversight, particularly in complex edge cases where		
490	regulatory guidance is silent or contradictory. Fur-		
491	thermore, the causal knowledge graph, while ex-		
492	tensive, relies on an initial "seed" of expert-labeled		
493	relationships, and expanding this pool to a broader		
494	range of practicing fire engineers would be neces-		
495	sary for wider industry validation.		
496	Ethical Considerations		
497	In safety-critical domains such as fire safety and		
498	evacuation planning, the automation of engineer-		
499	ing judgment carries significant ethical weight. Our		
500	framework is designed as a decision-support tool to		
501	enhance the reliability of simulation inputs, not as		
502	a replacement for certified professional engineers.		
503	The risk of over-reliance on AI-generated param-		
504	eters must be mitigated through clear transparency		
505	of sources, provided by our GraphRAG’s justifica-		
506	tion mechanism, ensuring that all derived values		
507	are traceable to authoritative engineering standards.		
508	Additionally, the use of open-weight models like		
509	Gemma 3 ensures that these tools can be deployed		
510	in secure, local environments, protecting sensitive		
511	building data from third-party exposure.		
512	References		
513	British Standards Institution. 2017. Bs 9999:2017		
514	fire safety in the design, management and use of		
515	buildings. code of practice . British Standard BS		
516	9999:2017, BSI, London, UK.		
517	Tivadar Danka and Peter Horvath. 2018. modal: A		
518	modular active learning framework for python . <i>arXiv</i>		
519	<i>preprint arXiv:1805.00979</i> .		
520	J. Dimyadi, M. Spearpoint, and R. Amor. 2008. Sharing		
521	Building Information using the IFC Data Model for		
522	FDS Fire Simulation . <i>Fire Safety Science</i> , 9:1329–		
523	1340.		
524	Dilan Durmus, Alberto Giretti, Ori Ashkenazi, Alessan-		
525	dro Carbonari, and Shabtai Isaac. 2024. The Role of		
	Large Language Models for Decision Support in Fire		526
	Safety Planning . Lille, France.		527
	Darren Edge, Ha Trinh, Newman Cheng, Joshua		528
	Bradley, Alex Chao, Apurva Mody, Steven Truitt,		529
	Dasha Metropolitansky, Robert Osazuwa Ness, and		530
	Jonathan Larson. 2024. From local to global: A		531
	graph rag approach to query-focused summarization .		532
	<i>arXiv preprint arXiv:2404.16130</i> .		533
	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,		534
	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang,		535
	and Haofen Wang. 2024. Retrieval-Augmented Gen-		536
	eration for Large Language Models: A Survey . <i>arXiv</i>		537
	<i>preprint</i> . ArXiv:2312.10997 [cs].		538
	Steven MV Gwynne, Erica Kuligowski, Michael Spear-		539
	point, and Enrico Ronchi. 2015. Bounding defaults		540
	in egress models . <i>Fire and Materials</i> , 39(4):335–352.		541
	Syed Muhammad Ali Hashmi, Tamas Fekete, and Hen-		542
	dro Wicaksono. 2024. Causal AI in the Automotive		543
	Industry: Impact Analysis Through Carbon Emission		544
	Case Study . In <i>2024 IEEE International Conference</i>		545
	<i>on Industrial Engineering and Engineering Manage-</i>		546
	<i>ment (IEEM)</i> , pages 1267–1271, Bangkok, Thailand.		547
	IEEE.		548
	Matthew Honnibal, Ines Montani, Sofie Van Lan-		549
	deghem, Adriane Boyd, and 1 others. 2020. spacy:		550
	Industrial-strength natural language processing in		551
	python .		552
	Haoyu Huang, Yongfeng Huang, Junjie Yang, Zhenyu		553
	Pan, Yongqiang Chen, Kaili Ma, Hongzhi Chen, and		554
	James Cheng. 2025. Retrieval-augmented genera-		555
	tion with hierarchical knowledge . <i>arXiv preprint</i>		556
	<i>arXiv:2503.10150</i> .		557
	Morgan J. Hurley, Daniel Gottuk, John R. Hall,		558
	Kazunori Harada, Erica Kuligowski, Milosh Pu-		559
	chovsky, José Torero, John M. Watts, and Christopher		560
	Wieczorek, editors. 2016. SFPE Handbook of Fire		561
	Protection Engineering . Springer, New York, NY.		562
	Emre Kiciman, Robert Ness, Amit Sharma, and Chen-		563
	hao Tan. 2023. Causal reasoning and large language		564
	models: Opening a new frontier for causality . <i>Trans-</i>		565
	<i>actions on Machine Learning Research</i> .		566
	Erica D. Kuligowski and Steve M. V. Gwynne. 2010.		567
	The Need for Behavioral Theory in Evacuation Mod-		568
	eling . In <i>Pedestrian and Evacuation Dynamics 2008</i> ,		569
	pages 721–732, Berlin, Heidelberg, Springer.		570
	Sihao Li, Jiali Wang, and Zhao Xu. 2024. Auto-		571
	mated compliance checking for BIM models based		572
	on Chinese-NLP and knowledge graph: an integrative		573
	conceptual framework . <i>Engineering, Construction</i>		574
	<i>and Architectural Management</i> , 32(6):3832–3856.		575
	Publisher: Emerald Publishing Limited.		576
	Jerry Liu. 2022. LlamaIndex .		577

578 Wensheng Lu, Keyu Chen, Ruizhi Qiao, and Xing Sun.
579 2025. Hichunk: Evaluating and enhancing retrieval-
580 augmented generation with hierarchical chunking.
581 *arXiv preprint arXiv:2509.11552*.

582 Linhao Luo, Zicheng Zhao, Gholamreza Haffari, Yuan-
583 Fang Li, Chen Gong, and Shirui Pan. 2024. Graph-
584 constrained reasoning: Faithful reasoning on knowl-
585 edge graphs with large language models. *arXiv*
586 *preprint arXiv:2410.13080*.

587 Soumya Madireddy, Lu Gao, Zia Ud Din, Kinam Kim,
588 Ahmed Senouci, Zhe Han, and Yunpeng Zhang. 2025.
589 [Large Language Model-Driven Code Compliance](#)
590 [Checking in Building Information Modeling](#). *Elec-*
591 *tronics*, 14(11):2146. Number: 11 Publisher: Multi-
592 disciplinary Digital Publishing Institute.

593 Neo4j, Inc. 2025. [Neo4j graph database](#).

594 Amir Rafe, Peter J. Lawrence, Ruggiero Lovreglio,
595 Michael Spearpoint, and Patrick A. Singleton. 2025.
596 [Enhancing Occupant Evacuation Simulation Using](#)
597 [LLMs and Retrieval-Augmented Generation](#). pages
598 389–400. Publisher: American Society of Civil En-
599 gineers.

600 Alexander Ratner, Stephen H. Bach, Henry Ehrenberg,
601 Jason Fries, Sen Wu, and Christopher Ré. 2017.
602 [Snorkel: rapid training data creation with weak su-](#)
603 [pervision](#). *Proc. VLDB Endow.*, 11(3):269–282.

604 Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh
605 Khanna, Anna Goldie, and Christopher D Manning.
606 2024. Raptor: Recursive abstractive processing for
607 tree-organized retrieval. In *The Twelfth International*
608 *Conference on Learning Representations*.

609 Jeffrey Tubbs, Gregory Harrington, and Yee Joe, editors.
610 2024. *NFPA 101: life safety code*. National Fire
611 Protection Association.

612 Upstage. 2024. Dp-bench: Document pars-
613 ing benchmark. [https://huggingface.co/](https://huggingface.co/datasets/upstage/dp-bench)
614 [datasets/upstage/dp-bench](https://huggingface.co/datasets/upstage/dp-bench). Accessed: 2025-
615 01-05.

616 Nengbo Wang, Xiaotian Han, Jagdip Singh, Jing Ma,
617 and Vipin Chaudhary. 2025. [CausalRAG: Integrating](#)
618 [Causal Graphs into Retrieval-Augmented Genera-](#)
619 [tion](#).

620 Xue Wu and Kostas Tsioutsoulis. 2024. Think-
621 ing with knowledge graphs: Enhancing llm rea-
622 soning through structured data. *arXiv preprint*
623 *arXiv:2412.10654*.

624 Appendices

625 A Prompt Engineering for Causal 626 Inference

627 To reduce hallucinations and ensure physical con-
628 sistency, we use a structured prompt that forces the
629 LLM to reason over retrieved causal graphs before
630 generating parameters.

```

631 Prompt Structure for Reasoning LLMs.
632
633 ## **LLM Query: Structured Reasoning-Based
634 Extraction for Walking Speed Distributions**
635
636 **[QUESTION]**
637
638 ### **Step 1: Consider These Causal
639 Relationships**
640 Before proceeding, ensure your reasoning aligns
641 with the following **causal relationships**:
642
643 [CAUSAL RELATIONSHIPS LIST]
644
645 ### **Step 2: Utilize Structured Knowledge
646 Extracted from Manuals**
647 The following **structured data** has been
648 extracted from fire safety manuals (NFPA,
649 SFPE, BS-9999, CIBS Guide E, and other
650 regulations):
651
652 [BASELINE VALUES]
653
654 ### **Step 3: Reason & Generate a Structured
655 Answer**
656 Using the provided **causal relationships** and
657 **structured extracted data**, determine the
658 most appropriate response.
659
660 - If numerical data is directly available,
661 provide **exact probability distributions**
662 (e.g., normal, log-normal, Weibull,
663 triangular, uniform, discrete) along with **
664 mean, standard deviation, and range**.
665
666 - If no direct numerical data is available,
667 analyze the causal mechanisms and **derive
668 reasonable estimates based on logical
669 constraints**.
670
671 - If conflicting values exist, provide the **
672 range of valid values** along with
673 justifications.
674
675 - If necessary, supplement reasoning with **peer-
676 reviewed studies**, but prioritize
677 structured data and causal reasoning first.
678
679 - If conflicting values exist between **manuals
680 & peer-reviewed studies**, explain the
681 differences and justify the most appropriate
682 choice.
683
684 - Do **not fabricate studies**--only reference
685 those explicitly available.
686
687 - Clearly indicate **assumptions made** if gaps
688 in the extracted knowledge exist.
689
690 ### **Step 4: Format Output as Follows**
691 1. **Final Answer** (Provide the most
692 appropriate statistical distribution &
693 values for walking speeds in different
694 environments)
695
696 2. **Justification & Reasoning** (Explain using
697 causal mechanisms & extracted knowledge)
698
699 3. **Comparison With Literature (If Available)**
700 (Compare with known fire safety studies)
701
702 4. **Causal Relationship Validation** (Ensure
703 output aligns with known causal dependencies
704 )

```

Listing 1: Structured Reasoning-Based Extraction Prompt

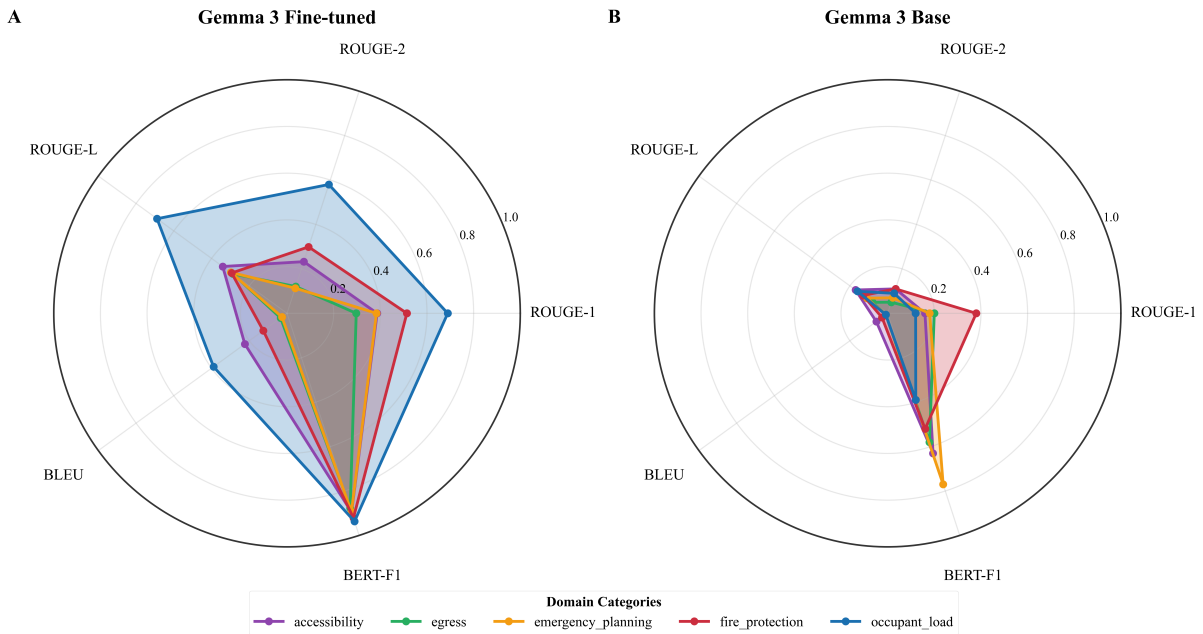


Figure 3: Radar charts comparing BERT-F1, ROUGE, and BLEU scores. (A) Gemma 3 Fine-Tuned shows balanced high performance across all metrics. (B) Gemma 3 Base exhibits significant degradation in specialized tasks like Occupant Load.

B Fine-Tuning Performance Analysis

Figure 3 presents a radar-based comparison of the fine-tuned Gemma 3 (4B) model (Panel A) against its base counterpart (Panel B) across five safety-critical domain categories: accessibility, egress, emergency planning, fire protection, and occupant load. Each radar axis corresponds to a complementary evaluation metric (ROUGE-1, ROUGE-2, ROUGE-L, BLEU, and BERT-F1), capturing both surface-level lexical overlap and deeper semantic alignment.

The fine-tuned model demonstrates substantial and consistent gains across all domains and metrics, with the most pronounced improvements observed in BERT-F1 and ROUGE-L, indicating enhanced semantic fidelity and clause-level understanding of regulatory and safety text. In particular, the occupant-load and fire-protection domains exhibit the largest performance envelopes, reflecting the effectiveness of domain-specific fine-tuning on numerically grounded and rule-intensive content. BLEU scores, while lower in absolute magnitude relative to semantic metrics, also increase uniformly after fine-tuning, suggesting improved local sequence consistency without overfitting to surface n-gram patterns.

In contrast, the base Gemma 3 model (Panel B) exhibits compressed radar profiles with limited sep-

aration across domains, especially for ROUGE-2, ROUGE-L, and BLEU. While moderate BERT-F1 scores indicate some baseline semantic capability, the overall narrow metric spread highlights the model’s limited ability to robustly generalize across heterogeneous safety standards without targeted adaptation.

C Network Visualizations

Figure 4 presents the Neo4j-based regulatory knowledge graph derived from NFPA standards, illustrating a representative subset comprising 500 nodes out of a total of 33,686 and 974 relationships out of 76,377. Figure 5 depicts the abstract causal dependency network used to constrain and guide reasoning about evacuation performance. Nodes correspond to interpretable simulation and behavioral variables (e.g., *Pre-Evacuation Time*, *Walking Speed*, *Visibility*, *Occupant Compliance*), while directed edges encode hypothesized or empirically supported causal influences. Central outcome variables such as *Evacuation Time* emerge as hubs, integrating effects from perceptual factors (smoke density, lighting), social dynamics (group cohesion, staff credibility), environmental cues, and individual characteristics (physical limitations, prior experience).

