

# BridgeEQA: Virtual Embodied Agents for Real Bridge Inspections

## Supplementary Material

### 1. Evaluating Image Citation Relevance for Human Alignment

To validate that our Image Citation Relevance metric aligns with human judgment, we conducted a manual annotation study on a randomly sampled set of 100 question-answer pairs from BridgeEQA. For each sample, we randomly perturbed the reference images by introducing varying numbers of random images from the original PDF document set and then randomly removing a varying number of images. This perturbation process generated image sets spanning the full relevance spectrum—from completely irrelevant to fully relevant to the question.

Three annotators independently labeled each sample on a 5-point scale, which we normalized to the 0.0-1.0 range to match the Image Citation Relevance output range. We then computed Image Citation Relevance scores for the same dataset using Gemini-2.5-flash as the evaluator model.

The Spearman correlation between the averaged human annotations and Image Citation Relevance scores was 0.817, demonstrating strong alignment between our automated metric and human judgment of image relevance.

### 2. Dataset Example

We provide an example from BridgeEQA in Figure 1. The reference images are images parsed from the source report and the condition rating is extracted from the answer. This particular scene graph for structure 0010 in Chelsea has a total of 53 nodes.

### 3. Dataset Creation Details

#### 3.1. Example Source Inspection Reports

We provide sample pages from Vermont bridge inspection report in Figure 2 as an example source report for BridgeEQA.

#### 3.2. Data Collection

We collected bridge inspection reports from the Vermont Agency of Transportation (VTTrans) public database, which contains unstructured PDF inspection reports covering bridges across Vermont. Each report documents the condition of a single unique bridge and includes inspector observations, condition ratings, and photographic documentation.

#### 3.3. Stage 1: Preprocess and Filter

The preprocessing stage applies quality control filters to ensure that selected reports contain sufficient visual documen-

tation for a meaningful infrastructure assessment.

**Report-Level Filtering.** Initial qualitative evaluation of the inspection reports revealed significant variability in visual documentation quality and comprehensiveness. To ensure sufficient visual coverage for meaningful condition assessment, we applied a minimum threshold of 20 images per report. Reports failing to meet this criterion were excluded based on several quality indicators:

- **Incomplete visual coverage:** Reports with fewer images often documented only limited perspectives of the bridge, missing critical structural components necessary for comprehensive assessment.
- **Low image quality:** Sparse image sets frequently exhibited poor resolution, unfavorable lighting conditions, or obstructed views that would hinder reliable condition evaluation.
- **Non-standard outlier conditions:** Some reports documented bridges that were demolished, under major reconstruction, or otherwise not representative of typical operational infrastructure.

**Page Filtering.** Page filtering removes the first two pages of each report, which typically contain administrative cover pages, title pages, and summary information without detailed inspection content or photographic documentation.

**Image-Level Filtering.** Within the quality-controlled reports, individual images underwent additional filtering. Images smaller than 200x200 pixels were systematically removed from the dataset. This threshold was established through empirical observation that sub-threshold images predominantly contained:

- Organizational logos and branding elements
- Document headers and administrative markings
- Thumbnails and preview images lacking structural detail

Such images provide minimal information for infrastructure condition assessment and could introduce noise into model training or evaluation.

**Random Sampling.** From the filtered pool of quality-controlled reports, we employed a random sampling strategy to select 200 reports for the final dataset. This sampling approach ensures representative coverage of Vermont’s bridge inventory while maintaining computational tractability for annotation and evaluation.

#### 3.4. Stage 2: Extract

The extraction stage processes filtered PDFs to obtain textual and visual content. Text extraction parses inspector notes, observations, and structured fields. Image extraction retrieves photographs meeting quality criteria, preserv-

## BridgeEQA Example: Structure 0010 in Chelsea

**Question:** Given the transverse cracking in the asphalt wearing surface where no joints are present, what are the resulting deterioration patterns—such as cracking and abrasion—observed on the abutments, and what is the overall condition rating for the substructure?

**Answer:** Asphalt over the structure has transverse cracking forming in various areas scattered throughout both sides over the abutments. The substructure below shows related deterioration. Reinforced concrete abutment #1 has multiple settlement cracks present, and several areas along the lower portions exhibit minor abrasion. Reinforced concrete abutment #2 has a full-height vertical crack, and its lower portions have heavy abrasion/scaling that penetrates as deep as approximately 15". Because all primary structural elements are sound but have minor section loss and cracking, the substructure is in **Fair Condition** with a rating of 5.

### Reference Images:

- 636045817baf4010d986b01a55baca1d.png
- 1de1058cd997112823da497df50cc66b.png
- 11a2bf060c346794067d43bc7718caad.png
- 3c435b2e6814c354430fee1480befc8d.png

**Extracted Condition Rating:** 5

### Scene Graph:

```
{
  "nodes": [
    {
      "image_name": "13b13f0ccab54674d56471f83082c47c.png",
      "central_focus": "Upstream view of Jail Brook from under the bridge.",
      "image_description": "View from under the bridge looking upstream...",
      "edges": [
        {
          "connected_to": "95b6875b4d3e1bcc4d524926a13b0157.png",
          "description_of_connection": "adjacent_to"
        },
        {
          "connected_to": "be36bcfe248d002bf1b6d2a62e6b20ae.png",
          "description_of_connection": "provides_context_for"
        }
      ]
    },
    ...
  ]
}
```

Figure 1. We provide an example from BridgeEQA on structure 0010 in Chelsea. This particular example has a scene graph with 53 nodes.

ing metadata about location and context. This stage yielded 9,586 images across 200 reports, averaging 47.93 per report.

### 3.5. Stage 3: Transform

The transformation stage structures the extracted content into standardized formats with ground truth annotations. We employ several vision-language models as zero-shot parsing tools to extract structured information from the inspection reports. Gemini 2.5 Flash and Gemini 2.5 Pro [3]

serves as the primary extraction model for its efficiency and quality, with no fine-tuning or training on the bridge inspection data. The models function purely as information extraction tools, parsing existing content rather than learning dataset-specific patterns. We found Gemini 2.5 Flash to frequently have parsing errors or hallucinations as context size's grew in this stage, as such we fall back to Gemini 2.5 pro to reprocess when these errors occur and drop reports if errors persist.

**Image Reference Mapping.** This component links pho-



Town 4 - ANDOVER  
District 2, 27 - WINDSOR County  
Owner: 3 - Town or Township Highway Agency  
Maintenance Responsibility: 3 - Town or Township Highway Agency

IDENTIFICATION	CLASSIFICATION
(1) State Name	00 - Vermont
(2) Bridge Number	10100000601
(3) Inventory Route	2 - District
(4) National Agency District	00 - Vermont
(5) Planning Agency District	00 - Vermont
(6) Plans Code	WILLIAMS RIVER
(7) Facility Category	0000
(8) Location	0.01 MI TO, 0.10 MI W
(9) Mile Post	0.00
(10) State Highway	00
(11) US Federal Route Number	00
(12) US State Route Number	00
(13) US County Route Number	00
(14) Local Route Number	00
(15) Inventory Status	0 - Not in Inventory
(16) Inventory Date	10/12/2023
(17) Inventory Type	0 - Routine
(18) Inventory Method	0 - Visual
(19) Inventory Frequency	0 - Annual
(20) Inventory Agency	00 - Vermont
(21) Inventory District	00 - Vermont
(22) Inventory County	00 - Windsor
(23) Inventory District	00 - District
(24) Inventory District	00 - District
(25) Inventory District	00 - District
(26) Inventory District	00 - District
(27) Inventory District	00 - District
(28) Inventory District	00 - District
(29) Inventory District	00 - District
(30) Inventory District	00 - District
(31) Inventory District	00 - District
(32) Inventory District	00 - District
(33) Inventory District	00 - District
(34) Inventory District	00 - District
(35) Inventory District	00 - District
(36) Inventory District	00 - District
(37) Inventory District	00 - District
(38) Inventory District	00 - District
(39) Inventory District	00 - District
(40) Inventory District	00 - District
(41) Inventory District	00 - District
(42) Inventory District	00 - District
(43) Inventory District	00 - District
(44) Inventory District	00 - District
(45) Inventory District	00 - District
(46) Inventory District	00 - District
(47) Inventory District	00 - District
(48) Inventory District	00 - District
(49) Inventory District	00 - District
(50) Inventory District	00 - District
(51) Inventory District	00 - District
(52) Inventory District	00 - District
(53) Inventory District	00 - District
(54) Inventory District	00 - District
(55) Inventory District	00 - District
(56) Inventory District	00 - District
(57) Inventory District	00 - District
(58) Inventory District	00 - District
(59) Inventory District	00 - District
(60) Inventory District	00 - District
(61) Inventory District	00 - District
(62) Inventory District	00 - District
(63) Inventory District	00 - District
(64) Inventory District	00 - District
(65) Inventory District	00 - District
(66) Inventory District	00 - District
(67) Inventory District	00 - District
(68) Inventory District	00 - District
(69) Inventory District	00 - District
(70) Inventory District	00 - District
(71) Inventory District	00 - District
(72) Inventory District	00 - District
(73) Inventory District	00 - District
(74) Inventory District	00 - District
(75) Inventory District	00 - District
(76) Inventory District	00 - District
(77) Inventory District	00 - District
(78) Inventory District	00 - District
(79) Inventory District	00 - District
(80) Inventory District	00 - District
(81) Inventory District	00 - District
(82) Inventory District	00 - District
(83) Inventory District	00 - District
(84) Inventory District	00 - District
(85) Inventory District	00 - District
(86) Inventory District	00 - District
(87) Inventory District	00 - District
(88) Inventory District	00 - District
(89) Inventory District	00 - District
(90) Inventory District	00 - District
(91) Inventory District	00 - District
(92) Inventory District	00 - District
(93) Inventory District	00 - District
(94) Inventory District	00 - District
(95) Inventory District	00 - District
(96) Inventory District	00 - District
(97) Inventory District	00 - District
(98) Inventory District	00 - District
(99) Inventory District	00 - District
(100) Inventory District	00 - District

ELEMENTS	DESCRIPTION	UNITS	TOTAL	CONDITION			
				CS1	CS2	CS3	CS4
12	Reinforced Concrete Deck	SF	1975	1240	600	135	0
1080	Delamination/Spall/Patched Area	SF	15	0	0	15	0
1120	Efflorescence/Rust Staining	SF	720	0	600	120	0
510	Wearing Surfaces	SF	1857	1299	0	558	0
3220	Crack (Wearing Surface)	SF	400	0	0	400	0
3230	Efflorescence (Wearing Surface)	SF	158	0	0	158	0
2000	Damage	LF	158	0	158	0	0
804	Concrete Facia	LF	158	33	100	25	0
1080	Delamination/Spall/Patched Area	LF	15	0	0	15	0
1120	Efflorescence/Rust Staining	LF	110	0	100	10	0

58 - Deck (S - SATISFACTORY CONDITION) - structural elements show some minor deterioration.) Reinforced concrete deck continues to deteriorate but remains in satisfactory condition. Bay #1 is fairly clean of defects. Bays #2 through #4 have small rust stains scattered throughout with larger areas of heavy saturation present. Deck has full depth concrete patches scattered throughout each bay. Downstream bays have some hairline cracks with efflorescence leakage present scattered throughout.

200 - Existing Wearing Surface Depth (5')

A21 - Deck Wearing Surface Condition (Poor) Asphalt is in poor condition. Heavy rutting cracking that has been sealed is present along the center line with various other cracking that has been sealed across the top surface is present. Multiple large patches are present along center line with small new potholes are forming that have exposed the membrane near the downhill side.

A24 - Deck Curb Condition (Good) Concrete curb is in fairly good condition having some minor wearing present. Fairly good amount of sediment and debris build up is present in front of both curbs with vegetation growth along the upstream curb.

A28 - Deck Drain Condition (Fair) Three (3) curb side drains are present along both the upstream and downstream curbs and are in fair condition. Fairly good amount of debris and sediment build up is present in front of all the drains blocking some of the deck drainage. Downstream curb side drains have heavy scaling below the drains along the fascia that has exposed the steel reinforcing.

A39 - Deck Facia Condition (Satisfactory) Concrete fascia is in okay condition. Downstream fascia has areas of heavy scaling that have exposed the steel reinforcing below the curb side drains with other scattered hairline to minor cracking with efflorescence leakage.

BC.05 Bridge Rating Condition Rating (SATISFACTORY - Widespread minor or isolated moderate defects.) Galvanized steel beam rail is in okay condition having areas of flattered out rail from pile up and scattered minor to moderate dents, scrapes and bends throughout. Reinforced concrete piles are in fairly good condition having some small chips of concrete missing along the edges and areas of light surface scaling and weathering. Post #5 and #10 along the upstream side and post #2 along the downstream side are missing connection bolts.

BC.08 Bridge Joints Condition Rating (NOT APPLICABLE - Bridge does not have deck joints.)



Deck & superstructure abutment 2



Abutment 2



Downstream



Upstream



Upstream of abutment 2



Upstream of abutment 1



Upstream rail from abutment 1



Deck & upstream rail from abutment 1

Maintenance Needs  
Date Reported: 09/22/2021  
Priority: 4 - Maintenance Finding - Next Status: Open  
Inspection Cycle: 1 - Annual  
Type of Work: 2 - General - Major rehabilitation project Component: General

Deficiency Description  
Structure is in need of a major rehabilitation project. Curb side along both the upstream and downstream sides should be cleaned and filled to prevent further corrosion to structure below. Asphalt over structure in allowing leakage to deck suffi below with multiple large patches and sealed cracks and deck should be reconditioned and paved with top surface of deck repairs. Approach and bridge rail have multiple areas of flattered out rail from pile up and should be considered for replacement. Beam #9 has areas along the entire side with heavy rust scaling and minor to moderate pitting from continuous leakage through deck curb drains and should be cleaned and painted to protect structural integrity. Structure could possibly need new concrete deck in near future.

Remarks



Wearing Surface on Abutment #2



Beam #5 near Midspan

Figure 2. Sample pages of BridgeEQA source report originating from the Vermont Agency of Transportation's (VTrans) inspection report for Structure 00006, located in Andover.

tographs to corresponding textual descriptions in inspector notes, supporting scene formation where multiple images document the same infrastructure component. This step is required to allow grounded questions that use real references for component names, such as Abutment 1.

**Condition Rating Extraction.** This component parses component-level NBI ratings from inspector assessments, providing ground truth labels on the standardized 0-9 scale [5].

**Inspector Note Preservation.** Inspector notes are preserved to maintain the original context and rationale for condition assessments, ensuring that ground truth annotations remain grounded in the source documentation. We

leverage these notes to ensure QA generation is grounded to real statements in the report.

### 3.6. Stage 4: Validate

Human quality control checks verify data integrity before QA generation. Additionally we test for any false or hallucinated image references. Parsing error detection identifies reports with corrupted text extraction, malformed condition ratings, or broken image-text mappings. When parsing errors or missing image references are detected, the report is automatically reprocessed using Gemini 2.5 Pro [3] as a fallback model for more robust extraction. Both Flash and Pro are used solely as zero-shot parsing tools without any

training or fine-tuning, ensuring that evaluation results reflect genuine visual reasoning capabilities rather than memorization. Reports that fail validation after reprocessing were removed from the dataset.

### 3.7. Stage 5: Generate QA

The final stage generates structured question-answer pairs for evaluation. Using a Gemini 2.5 Flash and Pro, we create questions grounded in the inspection report content, spanning condition assessment, component identification, and defect description tasks. Each answer includes the ground truth response sourced from inspector notes, references to supporting images, and the associated NBI condition rating when applicable. Quality checks verify that all referenced images exist and that answers are properly grounded in the available evidence.

### 3.8. Data Generation Validation

To ensure QA quality, we employ several evaluation metrics. We use the RAGAs [4] metrics: Faithfulness, which measures how well answers are grounded in the provided context, and Answer Relevancy, which assesses how effectively answers address the posed questions. We also incorporate the Answerability metric from RAGalyst [6] to determine whether questions can be adequately answered given the available context. To assess domain specificity, we employ LLM-as-a-Judge to determine Inspector Relevancy (0.0-1.0). This score measures the direct applicability of the question and its associated answer for bridge inspectors.

After evaluating all QA’s with Gemini-2.5-flash, we reach a Faithfulness of 0.997, an Answer Relevancy of 0.997, an Answerability of 0.996, and an Inspector Relevancy of 0.980. These high scores across all metrics demonstrate the overall high quality of the dataset.

### 3.9. Human Validation

To validate the automated filtering and processing pipeline, human evaluation was conducted on a random subsets of the processed reports. This manual inspection verified that the quality-controlled reports met the following standards:

- Sufficient visual coverage of critical bridge components
- Adequate image quality for condition assessment
- Consistency with typical operational bridge inspection documentation
- Accurate representation of the condition rating labels

## 4. Effects of Scene Graph Connectivity on Condition Rating

We provide the accuracy heatmap of each method and VLM by the number of nodes in the scene graph in Figure 3 and by the number of edges in the scene graph in Figure 4.

Table 1. Condition rating exact match accuracy (%) on open-source VLMs evaluated on BridgeEQA instances with fewer than 30 images.

Method	Qwen3-VL	Qwen3-VL	Nemotron-3
	8B-Thinking[1] <sup>†</sup>	30B-A3B[1]	Nano-30B-A3B[2]
Multi-Frame VLM	9.1	8.5	11.0
Socratic LLM w/ SG	36.4	6.1	11.0
EMVR VLM w/ SG Only	40.9	<b>29.3</b>	<b>23.2</b>

<sup>†</sup> Due to its limited context window, the Qwen3-VL 8B model was evaluated only on scenes with fewer than 30 images, a small fraction of the full dataset. These results are not directly comparable to other models reported in this paper.

Table 2. Condition rating within  $\pm 1$  accuracy (%) on open-source VLMs evaluated on BridgeEQA instances with fewer than 30 images.

Method	Qwen3-VL	Qwen3-VL	Nemotron-3
	8B-Thinking[1] <sup>†</sup>	30B-A3B[1]	Nano-30B-A3B[2]
Multi-Frame VLM	81.8	23.2	58.5
Socratic LLM w/ SG	72.7	30.5	48.8
EMVR VLM w/ SG Only	81.8	<b>76.8</b>	<b>70.7</b>

<sup>†</sup> Due to its limited context window, the Qwen3-VL 8B model was evaluated only on scenes with fewer than 30 images, a small fraction of the full dataset. These results are not directly comparable to other models reported in this paper.

Generally, the performance across methods decreases as the number of edges and nodes increase. This is due to the increased context sizes which is known to reduce VLM performance. However, EMVR performance degrades less at higher node and edge counts since EMVR mitigates the "lost in the middle" problem.

## 5. Open-Source Model Results

We extend our evaluation to open-source VLMs (Vision-Language Models) to assess generalizability beyond proprietary models. Given the large context windows required by our dataset, we omit Multi-Frame VLM w/ SG and EMVR VLM w/ Images + SG, as both require encoding images alongside the scene graph, which exceeds the context window of these models. Additionally, these models exhibited high failure rates in structured output generation, hallucinated function calls, and repeated the same actions in loops during agent execution. Due to these limitations, only a fraction of BridgeEQA could be tested. Results in 1 and 2 should therefore not be compared against the main paper results.

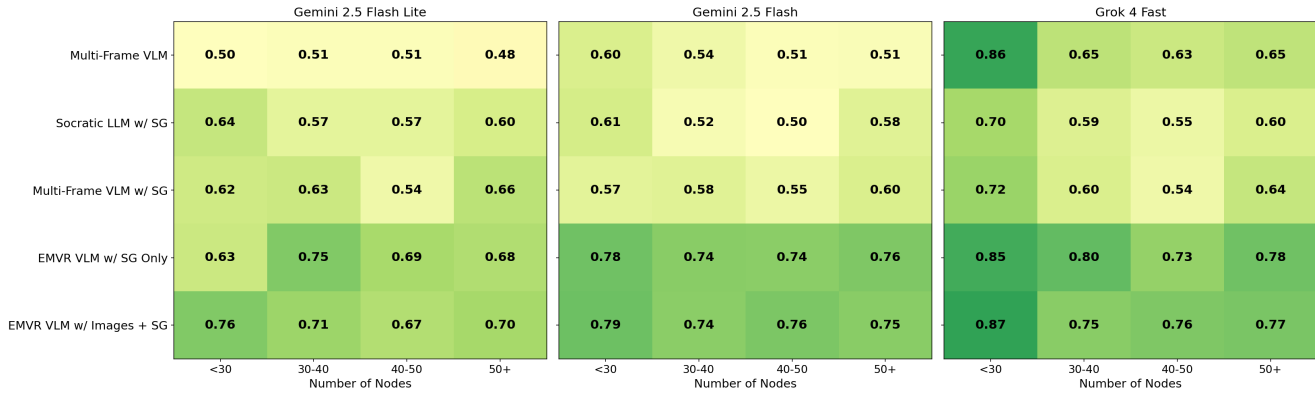


Figure 3. Condition rating within  $\pm 1$  accuracy heat map across method, VLM, and number of nodes.

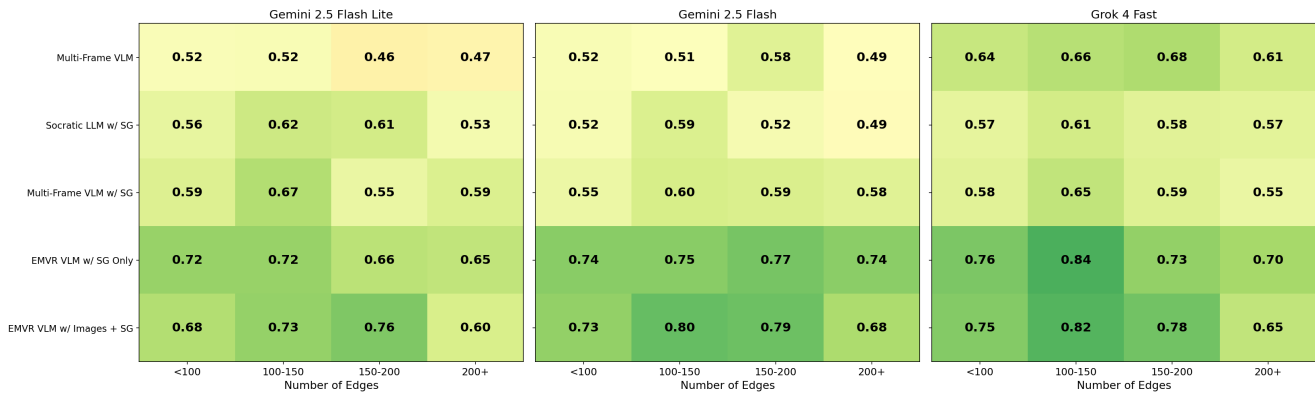


Figure 4. Condition rating within  $\pm 1$  accuracy heat map across method, VLM, and number of edges.

## References

- [1] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report, 2025. 4
- [2] Aaron Blakeman, Aaron Grattafiori, Aarti Basant, Abhibha Gupta, Abhinav Khattar, Adi Renduchintala, Aditya Vavre, Akanksha Shukla, Akhiad Bercovich, Aleksander Ficek, et al. Nvidia nemotron 3: Efficient and open intelligence. *arXiv preprint arXiv:2512.20856*, 2025. 4
- [3] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blis-  
tein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 2, 3
- [4] Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. RAGAs: Automated evaluation of retrieval augmented generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 150–158, St. Julians, Malta, 2024. Association for Computational Linguistics. 4
- [5] Federal Highway Administration. Recording and coding guide for the structure inventory and appraisal of the nation’s bridges. Technical Report FHWA-PD-96-001, U.S. Department of Transportation, Federal Highway Administration, 1995. 3
- [6] Joshua Gao, Quoc Huy Pham, Subin Varghese, Silwal Saurav, and Vedhus Hoskere. Ragalyst: Automated human-aligned agentic evaluation for domain-specific rag. *arXiv preprint arXiv:2511.04502*, 2025. 4