# Leveraging Multimodal LLMs for Building Condition Assessment from Street-View Imagery

Siyuan Yao $^{1[0000-0002-4093-193X]},$  Siavash Ghorbany $^{1[0000-0002-9588-0527]},$  Meghan Forstchen $^{1[0009-0003-7637-4092]},$  Alexis Korotasz $^{1[0009-0007-6560-068X]},$  Matthew Sisk $^{1[0000-0002-4141-9655]},$  Ming Hu $^{1[0000-0003-2583-1161]},$  and Chaoli Wang  $^{1[0000-0002-0859-3619]}$ 

University of Notre Dame, Notre Dame, IN 46556, USA {syao2,sghorban,mforstch,akorotas,msisk1,mhu1,chaoli.wang}@nd.edu

Abstract. We present a novel framework for automatically evaluating building conditions nationwide in the United States by leveraging large language models (LLMs) and Google Street View (GSV) imagery. By fine-tuning Gemma 3 27B on a modest human-labeled dataset, our approach achieves strong alignment with human mean opinion scores (MOS), outperforming even individual raters relative to the MOS benchmark in terms of SRCC and PLCC. To enhance efficiency, we apply knowledge distillation, transferring the capability of Gemma 3 27B to a smaller Gemma 3 4B model, which attains comparable performance with a  $3\times$  speedup. Further, we distill the knowledge into a CNN-based model (EfficientNetV2-M) and a transformer (SwinV2-B), delivering close performance while achieving a  $30\times$  speed gain. Our framework offers a flexible and efficient solution for large-scale building condition assessment, enabling high accuracy with minimal human labeling effort.

**Keywords:** Building condition evaluation  $\cdot$  Street-view imagery  $\cdot$  Computer vision  $\cdot$  Machine learning  $\cdot$  Multimodal large language models.

### 1 Introduction

The persistent shortage of affordable housing in the United States, aging infrastructure, and rising energy costs burden low-income households [2]. Retrofitting existing housing stock has emerged as a more feasible and cost-effective alternative to new construction, offering opportunities to improve thermal comfort, reduce utility bills, and mitigate health risks from extreme heat and cold [12]. Assessing the physical condition of buildings, especially the exterior envelope (i.e., façades), is critical to determining retrofit needs. However, audits remain resource-intensive, requiring manual inspections and data collection, limiting their practice at a nationwide scale.

Advances in computer vision [21, 6, 26, 7] have enabled automated analysis of street-view imagery, such as *Google Street View* (GSV) images, to identify *passive design indicators* (PDIs) like window-to-wall ratios, shading devices, and exterior material types. While these approaches leverage standard computer vision

models (e.g., ResNet) or vision-language models (e.g., BLIP) to classify exterior features with high accuracy, their capacity to simultaneously consider diverse visual elements indicating various features (e.g., paint, window, and structural conditions) and provide interpretable assessments remains limited. Large language models (LLMs), especially multimodal LLMs that integrate visual and textual reasoning, present a promising new frontier for interpretable and generalizable assessment. Aligning LLMs with human rater evaluations through fine-tuning on annotated street-view imagery could enable accurate and scalable assessments of building conditions, thereby supporting data-driven passive retrofit strategies across diverse urban contexts.

This paper presents a novel method for fine-tuning multimodal LLMs to automate exterior condition assessments for residential buildings, leading to a pipeline that fine-tunes multimodal LLMs to interpret building exteriors, enabling scalable, automated, and cost-effective assessments.

### 2 Related Work

Automated visual assessment of the built environment has become an increasingly active area of research, driven by the growing availability of street-view imagery and advances in machine learning. Traditional computer vision approaches, such as CNNs, have been employed to extract architectural features from images for estimating building condition [9,1,27]. However, these models often rely on narrow indicators, such as wall paint cracks, and struggle with low accuracy when integrating multiple factors due to limited generalizability.

More recently, vision-language models have demonstrated impressive zeroshot classification performance on visual tasks by leveraging cross-modal alignment between images and text. These models have been applied to identify building components or construction materials [26], as well as to estimate physical walkability [17] and perceived safety [25], directly from street-view imagery without needing extensive labeled data. However, their effectiveness in structured assessment tasks, such as condition scoring guided by formal criteria, remains limited due to insufficient capacity for multi-factor reasoning.

There is a growing interest in using LLMs to extract features of the built environment from street-view images. Several recent studies [4, 14, 19, 15] have employed ChatGPT to capture detailed building information, including external features and nearby environmental elements. However, ChatGPT's closed-source nature, lack of support for fine-tuning, and the high cost of API usage at scale limit its practicality for large datasets. Alternatively, fine-tuned open-source LLMs can be trained on labeled datasets for post-earthquake structural damage assessment, successfully performing tasks such as identifying damage severity and classifying affected components [13]. In this work, we fine-tune open-source LLMs for building condition assessment by leveraging diverse visual features that reflect natural aging and property upkeep. These include indicators like the state of windows, façade paint, and roof materials. Our system produces assessments

consistent with expert judgment by aligning model outputs with mean opinion scores (MOS) derived from multiple human raters.

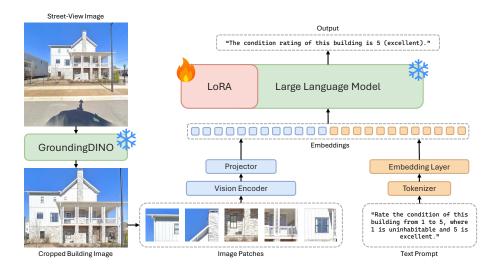


Fig. 1. Overview of our framework for evaluating building conditions.

## 3 Our Approach

We aim to evaluate the condition of a building from its GSV image. Experts in architectural research have established formal criteria for condition assessment using a five-point scale (higher ratings indicate better condition). Each rating is defined by detailed descriptors, covering various aspects such as the building's main structure, wall integrity, paint condition, roof, and window quality.

One possible approach is replicating human evaluation by training individual models for each component (e.g., roof, windows, façades), followed by a secondary algorithm to combine these assessments into an overall rating. However, this approach presents several limitations. First, many components may be partially occluded or absent from the image due to angle or visibility. Second, creating labeled datasets for each component would require significant manual effort from experts. Finally, integrating multiple separate models into a single scoring pipeline is complex and inefficient.

To address these challenges, we adopt a more direct strategy: evaluating the entire building holistically using a single model. This method assumes the model can effectively incorporate the full evaluation rubric and infer a corresponding rating. Following the method outlined in [26], and as illustrated in Figure 1, we preprocess images using GroundingDINO to isolate and crop individual buildings from street-view scenes. As found in the experiment in Section 4.2, we choose

#### 4 S. Yao et al.

Gemma 3 [23] as the multimodal LLM backbone to process each cropped image alongside a structured text prompt. In practice, in addition to the query, we include the formal rating criteria in the text prompt as follows:

- 1: Uninhabitable Likely unsuitable for rehabilitation; abandoned, fire-damaged, boarded-up, or vacant. Requires demolition.
- 2: Poor Requires substantial improvements, including major roof repairs, broken windows, bulging walls, or sagging foundations.
- 3: Adequate Requires basic cosmetic repairs, with no more than two issues such as painting/siding, trim, porch, minor roof improvements, or fence repair.
- 4: Good Structurally sound with good maintenance and no immediate repairs required. There may be no more than one minor issue, such as limited painting/siding replacement, minor porch repair/painting, or minor fence repair/painting.
- 5: Excellent Recently rehabilitated or remodeled; no repairs needed. New paint and roof in very good condition.

Formatting instructions for the expected output are also indicated in the text prompt, enabling the model to evaluate different aspects of the building, including paint, window, structure, and maintenance, according to the criteria, and then provide an overall numerical rating. Figure 2 shows five examples corresponding to ratings 1 to 5, respectively. This overall rating is directly compared with the MOS provided by human experts.



Fig. 2. Examples of buildings corresponding to each condition rating.

Within this framework, we employ two strategies to leverage Gemma 3: finetuning with expert supervision and knowledge distillation for efficiency.

First, to bring the model's predictions closer to the mean opinion of human experts, we fine-tune a strong base model using a small set of human-labeled data. Given the substantial size of modern LLMs and the computational constraints of our local hardware environment, full fine-tuning is not feasible.

To address this, we adopt a parameter-efficient fine-tuning (PEFT) [10] strategy, which selectively updates a small subset of the model's parameters while keeping the majority fixed. Specifically, we apply quantized low-rank adaptation (QLoRA) [5], a technique that enables fine-tuning using low-precision weights and low-rank adapters. This approach significantly reduces memory consumption and training cost, while preserving model performance.

Second, to enable efficient large-scale evaluation, we distill the capabilities of the largest Gemma 3 model (teacher) into smaller, more efficient models (students). For tasks requiring detailed text output, we train a smaller Gemma 3 to replicate the teacher's reasoning and formatting. We transfer knowledge to lightweight vision models for rating-only tasks for faster inference. Distillation is performed using pseudo-labels generated by the teacher on unlabeled building images, removing the need for additional human annotation while preserving alignment with expert judgments.

## 4 Experiments

#### 4.1 Dataset and Metrics

We collected 12,063 GSV images from six states: California, Florida, Georgia, Indiana, New York, and Texas. The images capture a wide range of building conditions and styles. We invited seven human raters to independently evaluate the building condition of 1,281 randomly selected images on a 1–5 scale. As a result, the labeled images received MOS ratings of 5 for 229 images, 4 for 582 images, 3 for 345 images, 2 for 105 images, and 1 for 20 images. The remaining 10,782 unlabeled images are used for knowledge distillation experiments. Though resolution varies due to cropping, all images are clear enough for reliable assessment, with a minimum resolution of  $600 \times 300$ .

Since the label distribution is imbalanced (e.g., 1's and 2's are scarce), we emphasize correlation-based evaluation, which is suitable for ordinal MOS labels and robust to skewed class frequencies. We employ the *Spearman's rank correlation coefficient* (SRCC) to evaluate the alignment between model-predicted ratings and human assessments. This non-parametric metric measures the monotonic relationship between two ranked variables, capturing how well the predicted ratings preserve the relative ordering of the ground truth ratings. SRCC is especially well-suited for subjective tasks, where the precise numerical rating may vary across raters, but the relative ranking remains meaningful. In equation,

$$SRCC = 1 - \frac{6\sum_{i=1}^{N} (x_i - y_i)^2}{N(N^2 - 1)},$$
(1)

where  $x_i$  and  $y_i$  denote the model-predicted rating and MOS from human ratings for the *i*-th image, respectively, and N denotes the number of images. A higher SRCC value indicates stronger agreement in ranking between the model and human raters, with  $\rho = 1$  indicating perfect rank correlation.

To assess the linear agreement between model predictions and human ratings, we also compute the *Pearson's linear correlation coefficient* (PLCC). Unlike SRCC, which measures monotonic rank alignment, PLCC quantifies the strength of a linear relationship between predicted ratings and ground truth ratings. It is defined as

PLCC = 
$$\frac{\sum_{i=1}^{N} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{N} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{N} (y_i - \bar{y})^2}},$$
 (2)

where  $x_i$  and  $y_i$  denote the model-predicted rating and MOS for the *i*-th image respectively, and  $\bar{x}$ ,  $\bar{y}$  are their corresponding means. N denotes the number of images. A PLCC value of 1 indicates perfect linear correlation, 0 indicates no correlation, and -1 indicates perfect inverse correlation. PLCC is useful when assessing whether the predicted ratings follow the correct ordering and approximate the correct magnitude.

#### 4.2 Zero-shot Evaluation of Multimodal LLMs

We compare the zero-shot performance of multimodal LLMs on building condition evaluation and assess their inference costs. We select the latest open-source models from the LLaVA [16], Mistral [20], LLaMA [24], Qwen [3], and Gemma [23] families. We evaluate two model sizes for LLaVA, Qwen, and Gemma to provide additional reference points. Experiments are run on a system equipped with four NVIDIA A40 GPUs (48 GB VRAM each). While most experiments are performed on a single GPU, models exceeding the memory capacity of one device in our system are run on two GPUs for inference. Note that the multi-GPU execution did not impact inference speed in our experiments. The test set comprises all 1,281 GSV images along with their MOS ratings. All models are given the same text prompt instructing them to evaluate paint, windows, structure, and maintenance, followed by an overall rating from 1 to 5 representing the building condition, which is compared with the MOS to calculate SRCC and PLCC.

As shown in Table 1, LLaMA 4 Scout 17B×16E ('B' denotes the number of parameters in *billions*, while 'E' stands for *experts* in a mixture of experts architecture) achieved the highest correlation with MOS, with an SRCC of 0.78 and a PLCC of 0.79, demonstrating strong agreement with human ratings, which has the same SRCC and PLCC as the average of human raters shown in Table 2. However, because the model requires more than 64 GB of memory, it had to be run on two GPUs in our setup. In comparison, Gemma 3 27B followed closely in performance, achieving a value of 0.77 for SRCC and PLCC. Notably, it requires only 18 GB of VRAM, making it compatible with a single commercial-grade GPU and therefore more practical for deployment on local machines with constrained computational resources. In this comparison, Gemma 3 27B shows the closest alignment with human ratings for building condition evaluation on a single GPU, outperforming larger models such as Mistral Small 3.2 24B, LLaVA 1.6 34B, Qwen 2.5 VL 32B, and Qwen 2.5 VL 72B.

**Table 1.** Comparison of multimodal LLMs. We report SRCC and PLCC with respect to the MOS, along with inference speed (tokens per second), GPU memory usage (VRAM, in gigabytes), and the number of GPUs (nGPU) used to run each model. Inference speed is expressed in tokens per second, ensuring consistency across models regardless of the length of generated text. The best values are highlighted in bold.

model	SRCC ↑	PLCC ↑	inference speed $\uparrow$	$\mathrm{VRAM}\downarrow$	nGPU↓
Mistral Small 3.2 24B	0.63	0.67	33.86	23.75	1
LLaMA 4 Scout 17B×16E	0.78	0.79	37.28	64.58	2
LLaVA 1.6 7B	0.41	0.41	88.46	5.52	1
LLaVA $1.6~34B$	0.61	0.61	25.83	25.92	1
Qwen~2.5~VL~32B	0.68	0.72	23.66	22.32	1
Qwen~2.5~VL~72B	0.73	0.76	11.02	49.77	2
Gemma 3 4B	0.45	0.46	97.37	4.87	1
Gemma 3 27B	0.77	0.77	26.64	18.48	1

**Table 2.** Comparison of individual human raters to the MOS, excluding each rater's own ratings to ensure fairness. We report SRCC and PLCC for each rater, as well as the average across all human raters.

rater	A	В	С	D	Е	F	G	average
SRCC ↑	0.78	0.80	0.80	0.76	$0.76 \\ 0.77$	0.74	0.81	0.78
PLCC ↑	0.80	0.80	0.81	0.78		0.76	0.82	0.79

Performance consistently declines among LLaVA, Qwen, and Gemma as model size is reduced. For Qwen 2.5 VL, the drop in both SRCC and PLCC from 72B to 32B is less than 0.1, which is much smaller than the 0.2 drop observed from 34B to 7B for LLaVA 1.6 or the drop of more than 0.3 from 27B to 4B for Gemma 3. As a result, although Gemma 3 4B demonstrated the highest efficiency in inference speed and the lowest VRAM usage, its predicted ratings deviate substantially from human ratings, making it unsuitable for use.

Overall, Gemma 3 27B is the most practical out-of-the-box open-source choice for building condition evaluation on a single GPU, combining near-human MOS alignment with moderate speed and GPU memory requirements.

**Table 3.** Comparison of different output formats. We report SRCC and PLCC, along with the average number of tokens in the response generated, the average response generation time (in seconds), and the prompt processing time per image (in seconds). The best values are highlighted in bold.

output	details & number	details & word	single number	single word
SRCC ↑	0.77	0.77	0.70	0.78
PLCC ↑	0.77	0.78	0.66	0.78
$\#$ response tokens $\downarrow$	87.82	79.89	2.00	2.00
response time $\downarrow$	3.29	3.00	0.08	0.08
processing time $\downarrow$	1.06	1.04	0.96	1.13

## 4.3 Optimization of Gemma 3

Building on the findings in Section 4.2, we investigate flexible ways to leverage Gemma 3 for building condition evaluation. We first examine whether different output formats in the prompt affect the performance of the Gemma 3 27B model. As shown in Table 3, we experiment with alternative formats for the overall rating, including using a word instead of a number, and restricting the model to output only a single number or word without additional text. This latter case is motivated by the fact that the number of tokens in the response directly impacts inference time. All experiments are evaluated on the same set of 1,281 images. The results in Table 3 indicate that when detailed descriptions are provided, the model's accuracy remains almost consistent regardless of whether the overall rating is expressed as a word or a number. When only the rating is needed, using a single word yields nearly the same accuracy but with much faster responses due to fewer generated tokens. However, accuracy drops significantly when restricted to outputting only a single number. We attribute this phenomenon to numbers with less semantic context than descriptive words, making it harder for the model to link to specific building conditions when used alone.

Next, we fine-tune Gemma 3 27B using MOS labels to improve its alignment with human ratings. Since the accuracy of the word-only response is comparable to that of detailed responses (see Table 3), we adopt the word-only format for these experiments to simplify loss computation based on MOS, which is translated from numerical ratings into their corresponding descriptive words. From the dataset of 1,281 images, the first 800 images are allocated for training and the remaining 481 for testing. As shown in Table 4, we vary the number of training images by randomly sampling subsets from the 800 training images. The slightly higher PLCC for the pre-trained Gemma 3 27B model compared to earlier results is attributable to the change in the test set. During fine-tuning, the model is quantized to 4-bit precision to reduce memory usage. Following the Gemma 3 official guidelines, the LoRA configuration uses a scaling factor of 16 and a dropout rate of 0.05 to prevent overfitting. The rank of the low-rank matrix is set to 16, and the warm-up ratio is set to 0.03 to improve training stability. LoRA limits training to roughly 16% of the model's parameters. The learning rate is set to  $5 \times 10^{-5}$ . The training batch size is 1, and only one epoch is run to avoid overfitting. The model is optimized with a next-token prediction loss, calculated as the cross-entropy between the predicted logits and the target labels.

**Table 4.** Comparison of different numbers of training images used to fine-tune the Gemma 3 27B model. We report SRCC and PLCC, along with the training time (in minutes). The best values are highlighted in bold.

# training images	0	100	200	300	400	500	600	700	800
SRCC ↑	0.78	0.77	0.78	0.80	0.81	0.82	0.82	0.83	0.83
PLCC ↑	0.79	0.78	0.79	0.81	0.81	0.83	0.81	0.82	0.82
training time $\downarrow$	_	4.52	8.20	12.68	16.57	21.48	25.38	30.52	34.68

Table 4 shows that fine-tuning Gemma 3 27B on 500 labeled images yields performance that surpasses the MOS alignment of all individual human raters (see Table 2). Increasing the number of training images beyond 500 does not yield definitive performance improvements, indicating that the model reaches a performance plateau. These results suggest that a fine-tuned Gemma 3 27B model trained on 500 images is sufficient to match or exceed human-level consistency, making it a practical replacement for manual rating in automated building condition evaluation.

## 4.4 Knowledge Distillation using Fine-tuned Gemma 3

We also experiment with response-based knowledge distillation using our Gemma 3 27B fine-tuned on 500 labeled images as a teacher model to fine-tune the Gemma 3 4B model, as well as several efficient vision models, including ResNet [8], MobileNetV3 [11], EfficientNetV2 [22], and Swin Transformer V2 [18]. For comparison, we prepare two training sets: one containing 10,782 building images labeled automatically by the fine-tuned Gemma 3 27B, and another containing 800 human-labeled images. Both are evaluated using the same 481 images. For the Gemma 3 4B model, we apply the same QLoRA fine-tuning procedure described in Section 4.3. For the vision models, we optimize using mean squared error (MSE) loss between the predicted ratings and the MOS, with a learning rate of  $1 \times 10^{-4}$  during fine-tuning. The vision models are trained for up to 10 epochs on the Gemma 3–labeled dataset and up to 100 epochs on the human-labeled dataset. We identify the epoch that achieves the highest SRCC on the test set and report its performance.

**Table 5.** Comparison of models fine-tuned on two training datasets labeled respectively by Gemma 3 27B and by human annotators. We report SRCC and PLCC for both datasets, along with the batch size used during training, inference speed (in images per second), and GPU memory usage at the inference stage (VRAM, in gigabytes). Inference speed is expressed in images per second because models other than Gemma 3 produce numeric outputs instead of text. The best values are highlighted in bold.

model	Gemma SRCC ↑	$\begin{array}{c} {\rm dataset} \\ {\rm PLCC} \uparrow \end{array}$	human SRCC ↑	dataset PLCC ↑	batch size	$\begin{array}{c} \text{inference} \\ \text{speed} \uparrow \end{array}$	VRAM↓
ResNet-50	0.68	0.66	0.52	0.53	32	69.41	1.87
MobileNetV3-L	0.65	0.66	0.45	0.46	32	69.91	1.72
EfficientNetV2-M	0.73	0.74	0.60	0.61	16	31.01	2.09
SwinV2-B	0.73	0.74	0.61	0.63	16	32.61	2.36
Gemma 3 4B	0.81	0.80	0.74	0.73	1	3.05	9.10

As shown in Table 5, knowledge distillation effectively transfers the capabilities of the fine-tuned Gemma 3 27B model to smaller models. For all models, fine-tuning on the automatically labeled dataset leads to consistently higher performance than fine-tuning on the human-labeled dataset. The approach performs

well across different architectures, including CNNs (ResNet-50, MobileNetV3-L, EfficientNetV2-M), transformers (SwinV2-B), and multimodal LLMs (Gemma 3 4B). This result indicates that the fine-tuned Gemma 3 27B model can serve as a dependable substitute for human annotators in creating a large-scale dataset for building condition evaluation, helping to overcome the scarcity of annotated data while significantly reducing labeling costs. Notably, EfficientNetV2-M and SwinV2-B achieve SRCC and PLCC values above 0.7, comparable to the base Gemma 3 27B model, while delivering more than 30× faster inference. Finetuning Gemma 3 4B on the automatically labeled dataset enables it to achieve SRCC and PLCC values exceeding the base Gemma 3 27B model, with the added benefit of roughly 3× faster inference. In this experiment, the fine-tuned LoRA adapters were not merged into the pre-trained model prior to inference. Instead, they were loaded separately at runtime, allowing us to keep only the compact LoRA files rather than full model checkpoints. This approach increased VRAM usage from less than 5 GB to over 9 GB, as both the base model and the adapter were held in memory during inference. Still, users can merge the adapters into the base model, reducing the VRAM consumption to a level close to the pretrained model. Leveraging knowledge distillation from the fine-tuned Gemma 3 27B model enables scalable, fully automated evaluation with the flexibility to select models for different performance-efficiency trade-offs, making it practical to process datasets containing millions of images within reasonable timeframes.

## 5 Conclusions and Future Work

We present a novel framework for automated building condition evaluation from GSV images across the United States using multimodal LLMs with minimal human annotation. We benchmark multiple leading open-source multimodal LLMs, identify the most effective model for aligning with expert ratings, and explore techniques, such as prompt refinement, targeted fine-tuning, and knowledge distillation, to enhance reliability and significantly improve efficiency, making million-scale dataset evaluation practical. The framework provides a flexible selection of methods, allowing users to balance performance and efficiency according to task-specific priorities. While this approach is the first to evaluate building quality at a large scale using LLMs, a promising future direction is to explore their application to broader aspects of buildings.

Still, our methods have a few limitations. First, the MOS ratings are derived from a specific group of skilled and novice raters, which may not fully capture broader subjective consensus or ensure complete fairness. While this may be sufficient for meeting the needs of a small target group, broader application scenarios require input from a larger and more diverse rater pool. Second, output-based knowledge distillation demands large quantities of raw images to generate pseudo-labeled image—rating pairs from real-world data. While feature-based knowledge distillation may offer greater efficiency, differences in network architectures introduce challenges in maintaining compatibility between the teacher and student feature representations. We consider exploring this as

a future research direction. Third, given that LLM performance is highly dependent on prompt construction, future work will explore more sophisticated prompt engineering strategies beyond the plain-language criteria, query, and format descriptions used in this study. Finally, future work could examine potential biases in LLM outputs stemming from various image attributes, both low-level factors such as brightness and contrast, and high-level characteristics such as architectural style or building type, to better understand influences on model predictions.

**Acknowledgments.** This research was supported in part by the U.S. National Science Foundation through grants IIS-1955395, IIS-2101696, OAC-2104158, IIS-2401144, and CNS-2430623, the University of Notre Dame's Just Transformations to Sustainability Initiative, Lucy Family Institute for Data & Society Health Equity Data Lab, and Notre Dame-IBM Technology Ethics Lab.

## References

- Amrouni Hosseini, M., Ravanshadnia, M., Rahimzadegan, M., Ramezani, S.: Next-generation building condition assessment: BIM and neural network integration. Journal of Performance of Constructed Facilities 38(6), 04024050 (2024)
- 2. Anacker, K.B.: Introduction: Housing affordability and affordable housing. International Journal of Housing Policy 19(1), 1–16 (2019)
- 3. Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al.: Qwen technical report. arXiv preprint arXiv:2309.16609 (2023)
- Cheng, Y., Yin, Z., Li, D., Li, Z.: Assessing urban safety: A digital twin approach using streetview and large language models. In: Proceedings of IEEE Vehicular Technology Conference. pp. 1–5 (2024)
- Dettmers, T., Pagnoni, A., Holtzman, A., Zettlemoyer, L.: QLoRA: Efficient finetuning of quantized LLMs. In: Proceedings of Advances in Neural Information Processing Systems. pp. 10088–10115 (2023)
- Ghorbany, S., Hu, M., Yao, S., Wang, C., Nguyen, Q.C., Yue, X., Alirezaei, M., Tasdizen, T., Sisk, M.: Examining the role of passive design indicators in energy burden reduction: Insights from a machine learning and deep learning approach. Building and Environment 250, 111126 (2024)
- Ghorbany, S., Hu, M., Yao, S., Sisk, M., Wang, C., Zhang, K., Nguyen, Q.C.: Data driven assessment of built environment impacts on urban health across United States cities. Scientific Reports 15, 19998 (2025)
- 8. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)
- 9. Hoang, N.D.: Image processing-based recognition of wall defects using machine learning approaches and steerable filters. Computational Intelligence and Neuroscience **2018**, 7913952 (2018)
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., Gelly, S.: Parameter-efficient transfer learning for NLP. In: Proceedings of IEEE International Conference on Machine Learning. pp. 2790–2799 (2019)

- Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al.: Searching for MobileNetV3. In: Proceedings of IEEE International Conference on Computer Vision. pp. 1314–1324 (2019)
- 12. Hu, M., Ghorbany, S., Yao, S., Wang, C., Sisk, M.: BUILT2AFFORD: Machine-learning-driven passive retrofits for affordable housing. In: Proceedings of Architectural Research Centers Consortium Annual Conference (2025)
- 13. Jiang, Y., Wang, J., Shen, X., Dai, K.: Large language model for post-earthquake structural damage assessment of buildings. Computer-Aided Civil and Infrastructure Engineering (2025)
- 14. Li, Z., Su, Y., Wang, H., Zhao, W.: BuildingView: Constructing urban building exteriors databases with street view imagery and multimodal large language mode. arXiv preprint arXiv:2409.19527 (2024)
- Liang, X., Xie, J., Zhao, T., Stouffs, R., Biljecki, F.: OpenFACADES: An open framework for architectural caption and attribute data enrichment via street view imagery. arXiv preprint arXiv:2504.02866 (2025)
- 16. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: Proceedings of Advances in Neural Information Processing Systems (2023)
- 17. Liu, X., Haworth, J., Wang, M.: A new approach to assessing perceived walkability: Combining street view imagery with multimodal contrastive learning model. In: Proceedings of ACM SIGSPATIAL International Workshop on Spatial Big Data and AI for Industrial Applications. pp. 16–21 (2023)
- Liu, Z., Hu, H., Lin, Y., Yao, Z., Xie, Z., Wei, Y., Ning, J., Cao, Y., Zhang, Z., Dong, L., et al.: Swin Transformer V2: Scaling up capacity and resolution. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. pp. 12009–12019 (2022)
- Malekzadeh, M., Willberg, E., Torkko, J., Toivonen, T.: Urban attractiveness according to ChatGPT: Contrasting AI and human insights. Computers, Environment and Urban Systems 117, 102243 (2025)
- 20. Mistral AI Team: Mistral Small 3 (2025), https://mistral.ai/news/mistral-small-3
- 21. Starzyńska-Grześ, M.B., Roussel, R., Jacoby, S., Asadipour, A.: Computer vision-based analysis of buildings and built environments: A systematic review of current approaches. ACM Computing Survey 55(13s), 284:1–284:25 (2023)
- Tan, M., Le, Q.: EfficientNetV2: Smaller models and faster training. In: Proceedings of IEEE International Conference on Machine Learning. pp. 10096–10106 (2021)
- Team, G., Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M.S., Love, J., et al.: Gemma: Open models based on Gemini research and technology. arXiv preprint arXiv:2403.08295 (2024)
- 24. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: LLaMA: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)
- 25. Wang, X., Gilvear, A., Li, Y., Ilyankou, I.: Can CLIP see safe streets? comparing human and VLM perceptions of walkability and safety. In: Proceedings of AGILE Walking the X-min City Workshop (2025)
- Yao, S., Ghorbany, S., Sisk, M., Hu, M., Wang, C.: Leveraging zero-shot learning on street-view imagery for built environment variable analysis. In: Proceedings of International Symposium on Visual Computing. pp. 243–254 (2024)
- Zou, S., Wang, L.: Detecting individual abandoned houses from Google street view:
   A hierarchical deep learning approach. ISPRS Journal of Photogrammetry and Remote Sensing 175, 298–310 (2021)