

---

# Assessing Building Heat Resilience Using UAV and Street-View Imagery with Coupled Global Context Vision Transformer

---

**Steffen Knoblauch**

Heidelberg University

steffen.knoblauch@uni-heidelberg.de

**Ram Kumar Muthusamy**

Heidelberg University

ramkumar.muthusamy@uni-heidelberg.de

**Hao Li**

National University of Singapore

hao.li@nus.edu.sg

**Iddy Chazua**

OpenMap Development Tanzania

iddy.chazua@omdtz.or.tz

**Benedcto Adamu**

OpenMap Development Tanzania

benedcto.adamu@omdtz.or.tz

**Innocent Maholi**

OpenMap Development Tanzania

innocent.maholi@omdtz.or.tz

**Alexander Zipf**

Heidelberg University

zipf@uni-heidelberg.de

## Abstract

Climate change is intensifying human heat exposure, particularly in densely built urban centers of the Global South. Low-cost construction materials and high thermal-mass surfaces further exacerbate this risk. Yet scalable methods for assessing such heat-relevant building attributes remain scarce. We propose a machine learning framework that fuses openly available unmanned aerial vehicle (UAV) and street-view (SV) imagery via a coupled global context vision transformer (CGCViT) to learn heat-relevant representations of urban structures. Thermal infrared (TIR) measurements from HotSat-1 are used to quantify the relationship between building attributes and heat-associated health risks. Our dual-modality cross-view learning approach outperforms the best single-modality models by up to 9.3%, demonstrating that UAV and SV imagery provide valuable complementary perspectives on urban structures. The presence of vegetation surrounding buildings (versus no vegetation), brighter roofing (versus darker roofing), and roofing made of concrete, clay, or wood (versus metal or tarpaulin) are all significantly associated with lower HotSat-1 TIR values. Deployed across the city of Dar es Salaam, Tanzania, the proposed framework illustrates how household-level inequalities in heat exposure—often linked to socio-economic disadvantage and reflected in building materials—can be identified and addressed using machine learning. Our results point to the critical role of localized, data-driven risk assessment in shaping climate adaptation strategies that deliver equitable outcomes.

## 1 Introduction

Climate change is driving more frequent and intense heat waves, posing a growing public health threat [1, 2, 3]. Urban heat islands, formed by human-built environments with extensive impervious

surfaces, further exacerbate these risks [4]. Although heat exposure occurs both outdoors and indoors, the majority of daily heat burden is experienced indoors [5, 6]. Indoor heat resilience is strongly influenced by socio-economic factors, particularly in the Global South, where building insulation standards are often absent or inadequate. Consequently, socio-economically disadvantaged households disproportionately bear the health impacts of heat exposure despite contributing minimally to the underlying drivers of climate change [7, 8]. Designing effective and equitable heat resilience strategies thus requires robust methods to monitor and characterize technical building attributes at scale.

Satellite imagery offers a promising avenue for large-scale extraction of technical building characteristics [9]. While global building footprint datasets derived from satellite data provide extensive geometric information, they typically lack detailed attributes beyond basic shape and size. OpenStreetMap enables manual tagging of individual buildings, but its coverage and consistency are limited by the reliance on volunteer contributions. Recent studies highlight the value of SV imagery in capturing heat-relevant building features not visible from above, such as wall materials [10, 11, 12, 13, 14], colors [15, 16], number of floors [17], and vegetation cover [18, 19]. Complementing this, UAV imagery—with its higher spatial resolution—outperforms satellite data in detailing roofing materials, often obscured in SV imagery, especially within dense urban canyons [20, 21]. Open data platforms like OpenAerialMap and Panoramax facilitate the sharing and integration of such diverse UAV and SV datasets, advancing scalable mapping of technical building attributes. However, the combined use of these complementary data sources remains underexplored [22, 23]. To date, we are not aware of any study that has applied dual-modality learning approaches [24] integrating aerial and SV imagery to extract building-specific features relevant to heat resilience. However, filling this research gap is essential for developing targeted and equitable climate mitigation strategies at the local scale.

## 2 Materials and methods

This study presents the development of a machine learning framework to extract heat-relevant building attributes from openly accessible geospatial datasets (cf. Fig. 1). Inputs included (i) SV panoramic imagery, (ii) building footprints, and (iii) high-resolution UAV imagery to construct a cross-view representation of buildings. A CGCViT was trained to classify buildings by structural openness, number of floors, vegetation, wall material, and roofing material. These attributes, together with distance to surrounding buildings and roof and wall brightness, were statistically associated with HotSat-1 TIR values to identify priority targets for building-level heat mitigation by revealing features most strongly linked to lower thermal exposure. The framework was applied to the Msimbazi River Delta, Dar es Salaam, Tanzania—a densely populated floodplain (530,837 inhabitants) featuring compact urban development and diverse building typologies, providing an ideal context for investigating inequalities in urban heat exposure.

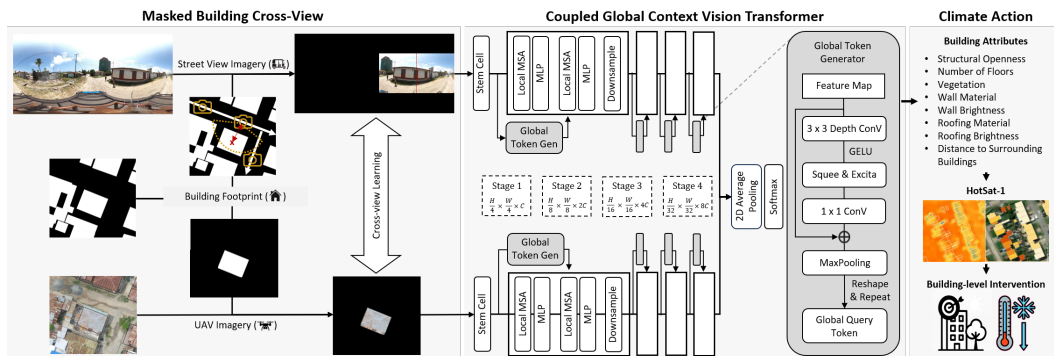


Figure 1: Overview of the proposed CGCViT framework for extracting heat-relevant building attributes from paired UAV and SV imagery. Building footprints guide the extraction of cross-view image pairs, which are processed in parallel GCvIT branches with global token integration. Predicted attributes include structural openness, number of floors, vegetation, wall/roofing material and brightness, and distance to surrounding buildings. These are statistically associated with HotSat-1 TIR values to inform heat mitigation strategies at building-level.

## 2.1 Masked building cross-view

SV imagery were collected via the Panoramax API in collaboration with OpenMap Development Tanzania (OMDTZ). A GoPro Max 360° camera mounted on a tricycle captured imagery at  $2048 \times 2048$  pixels at sub-4-meter intervals along accessible roads (cf. Fig. 4). Data acquisition occurred in two phases: Oct 28–Nov 12, 2024, and Jan 6–Feb 8, 2025 [25]. Coverage was limited by inaccessible narrow footpaths in informal settlements and restricted institutional sites; these were excluded (cf. Fig. 5a). UAV imagery, sourced from OpenAerialMap and collected jointly with OMDTZ, was acquired using a DJI Mavic 2 Pro drone flying at 150 m altitude during a two-week campaign starting October 25, 2023 (cf. Fig. 3). High-accuracy GPS ground control points ensured geospatial precision, yielding a 19.2 km<sup>2</sup> orthomosaic with 2.4 m horizontal and 1.4 m vertical accuracy and an average ground sampling distance of 9 cm. Building footprints from Geofabrik (downloaded April 8, 2025) were used to spatially align aerial and SV data across 63,844 buildings, of which 42,135 (65.98 %) were residential. Our analysis focused exclusively on residential structures, given their critical role in shaping household-level heat resilience. To ensure matching between aerial and SV imagery, the dataset was restricted to residential buildings with centroids within 30 m of a SV capture point and unobstructed frontal views, identified via a nearest-neighbor algorithm using building footprints. This filtering yielded 4,965 buildings with usable top- and front-view imagery (cf. Fig. 5b). To improve model focus, aerial and SV imagery were masked using building footprints. For SV imagery, driving direction metadata localized each building’s angular segment within the 360° image, enabling precise façade masking.

## 2.2 Coupled global context vision transformer

We trained the CGCViT on 2,000 masked cross-view building image pairs ( $256 \times 256$  px) with manual annotations for five classification tasks: structural openness, number of floors, surrounding greenery, wall material, and roofing material (cf. Fig. 6, Tab. 1). Wall and roof brightness were computed from mean RGB values of masked regions, and building density was estimated from mean distance to the four nearest OSM footprints (cf. Fig. 7, Fig. 8). The dataset was split into 70%/15%/15% train/val/test, with data augmentation applied to minority classes. CGCViT processes UAV and SV imagery in parallel GCViT-Tiny [26] streams, each with four stages of local and global self-attention [27, 28], capturing fine-grained patterns and long-range dependencies (cf. Fig. 1). A global query token is injected into local attention to aggregate context across regions via

$$\mathbf{G} = \text{Softmax}\left(\frac{\mathbf{g}_q \mathbf{k}^\top}{\sqrt{s}} + \mathbf{p}\right) \mathbf{v}, \quad (1)$$

where  $\mathbf{g}_q$  is the global query,  $\mathbf{k}$  and  $\mathbf{v}$  are key and value matrices,  $s$  is a scaling factor, and  $\mathbf{p}$  is a learnable relative positional embedding. Features from both streams are concatenated, pooled, and classified using softmax cross-entropy

$$\mathcal{L}_{\text{softmax}} = - \sum_{c=1}^C y_c \log \frac{e^{z_c}}{\sum_{j=1}^C e^{z_j}}. \quad (2)$$

For imbalanced tasks, focal loss was applied. Stratified 5-fold spatial cross-validation prevents geographic leakage. Model performance (support weighted-F1) was compared across UAV-only, street-view-only, and fused multi-modal cross-view learning settings. We assessed the statistical relationships between predicted building attributes and HotSat-1 TIR values (3.5 m spatial resolution; acquired June–December 2023) using the Kruskal–Wallis test for categorical variables and Pearson’s correlation coefficient ( $r$ ) for numerical variables.

## 3 Results and discussion

Cross-view learning using CGCViT consistently improved classification performance over the best single-modality models (UAV-only or SV-only), with gains in support weighted-F1 score ranging from 0% to 9.3% depending on the attribute. The most notable improvement occurred for vegetation classification (+9.3%), followed by structural openness (+7.7%) and number of floors (+7.1%). Roofing material classification showed a modest gain (+3.7%), while wall material classification exhibited no improvement over the best single-modality model (cf. Table 2). Vegetation classification

benefited substantially from the multi-modal approach, as UAV imagery captures vegetation in backyards or on rooftops that are often out of reach for SV; conversely, SV imagery provides valuable views of vegetation located underneath shelters or on window ledges, complementing the UAV perspective. Number of floors was best predicted from SV imagery alone; however, incorporating UAV data via cross-view fusion yielded additional gains, likely because larger building footprints visible in UAV imagery tend to be associated with greater vertical building extent. Roofing material classification benefited markedly from UAV imagery compared to SV alone, but the combined multi-modal model achieved the highest accuracy. The improvement over the UAV-only model presumably reflects the predominance of low-rise buildings and wide streets in our study area, which make roofing details more discernible in SV imagery and thus provide complementary information to the UAV perspective—a condition that may not generalize to more complex urban environments. Wall material classification remains challenging, with lower accuracy overall, probably due to the near-complete invisibility of walls in UAV imagery and high rates of building obstruction in SV imagery (e.g., by vehicles, fences, or vegetation), both of which constrain the discriminative information available to single-modality models. For structural openness, modality-specific differences were minimal; nevertheless, the multi-modal model achieved modest F1-score improvements despite the inherent difficulty of the task, compounded by ambiguous class boundaries identified during annotation. Beyond classification accuracy, our analysis revealed significant associations between certain predicted building attributes and HotSat-1 TIR values (Fig. 2). In particular, buildings with surrounding vegetation, roofing materials such as concrete, clay, or wood, and higher roof brightness exhibited lower mean TIR values, indicating a strong building-level cooling effect.

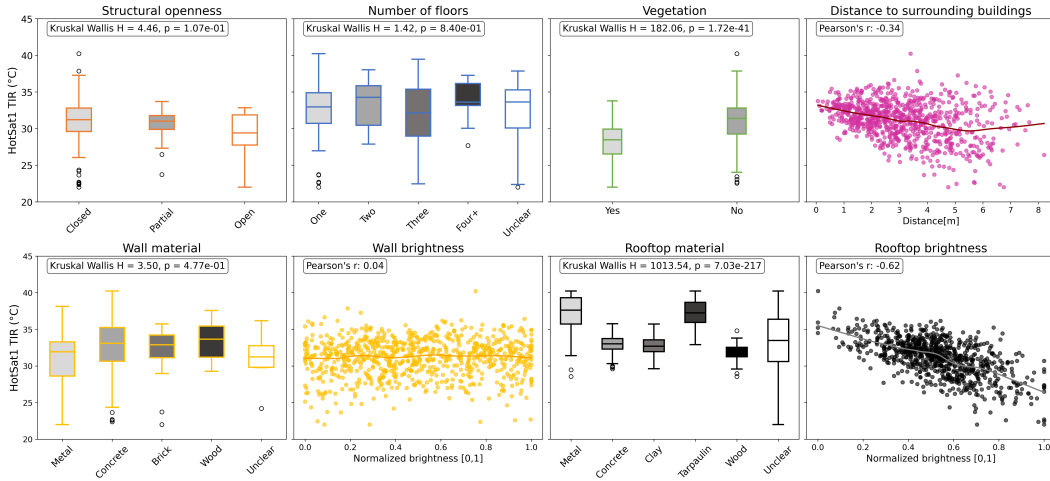


Figure 2: Boxplots and scatter plots illustrate the relationships between HotSat-1 TIR values and building attributes identified in this study. They highlight a significant drop in mean Hotsat-1 TIR values for buildings surrounded by vegetation, as well as the important role of roofing materials and brightness.

## 4 Conclusion

UAV and SV imagery, while not fully orthogonal, provide complementary perspectives that—when integrated via cross-view learning—improve the identification of heat-relevant building attributes. We find that greater vegetation cover, roofing materials such as concrete, clay, or wood, and higher roof brightness are consistently associated with lower HotSat-1 TIR values. This suggests that low-cost, household-level interventions—such as applying reflective roof coatings, providing seeds and water for small urban gardens, or supporting upgrades of heat-absorbing materials—can effectively reduce heat exposure. These scalable, machine-learning-derived insights enable targeted climate action while promoting equitable outcomes for households most vulnerable to climate-related impacts.



## References

- [1] Matthew Abunyewah, Thayaparan Gajendran, Michael Odei Erdiaw-Kwasie, Charles Baah, Seth Asare Okyere, and Amila Kasun Sampath Udage Kankanamge. The multidimensional impacts of heatwaves on human ecosystems: A systematic literature review and future research direction. *Environmental Science & Policy*, 165:104024, 2025.
- [2] Samuel Lüthi, Christopher Fairless, Erich M. Fischer, Noah Scovronick, Armstrong Ben, Micheline De Sousa Zanotti Stagliorio Coelho, Yue Leon Guo, Yuming Guo, Yasushi Honda, Veronika Huber, Jan Kysely, Eric Lavigne, Dominic Royé, Niilo Ryti, Susana Silva, Aleš Urban, Antonio Gasparrini, David N. Bresch, and Ana M. Vicedo-Cabrera. Rapid increase in the risk of heat-related mortality. *Nature communications*, 14(1):4894, 2023.
- [3] Kristie L. Ebi, Anthony Capon, Peter Berry, Carolyn Broderick, Richard de Dear, George Havenith, Yasushi Honda, R. Sari Kovats, Wei Ma, Arunima Malik, Nathan B. Morris, Lars Nybo, Sonia I. Seneviratne, Jennifer Vanos, and Ollie Jay. Hot weather and heat extremes: health risks. *Lancet (London, England)*, 398(10301):698–708, 2021.
- [4] Shengjun Gao, Yunhao Chen, Deliang Chen, Bin He, Adu Gong, Peng Hou, Kangning Li, and Ying Cui. Urbanization-induced warming amplifies population exposure to compound heatwaves but narrows exposure inequality between global north and south cities. *npj Climate and Atmospheric Science*, 7(1), 2024.
- [5] Jalonnie L. White-Newsome, Brisa N. Sánchez, Olivier Jolliet, Zhenzhen Zhang, Edith A. Parker, J. Timothy Dvonch, and Marie S. O’Neill. Climate change and health: indoor heat exposure in vulnerable populations. *Environmental research*, 112:20–27, 2011.
- [6] Ebenezer F. Amankwaa, Ben M. Roberts, Peter Mensah, and Katherine V. Gough. Impact of roofing materials on school temperatures in tropical africa. *Buildings and Cities*, 6(1):139–157, 2025.
- [7] Andrew Norton and Robin Mearns. Social dimensions of climate change : Equity and vulnerability in a warming world. new frontiers of social policy, 2010.
- [8] Yuan Yuan, Mattheos Santamouris, Dong Xu, Xiaolei Geng, Chengwei Li, Wanqing Cheng, Ling Su, Peng Xiong, Zhengqiu Fan, Xiangrong Wang, and Chuan Liao. Surface urban heat island effects intensify more rapidly in lower income countries. *npj Urban Sustainability*, 5(1), 2025.
- [9] Xiao Xiang Zhu, Sining Chen, Fahong Zhang, Yilei Shi, and Yuanyuan Wang. Globalbuilding-atlas: An open global and complete dataset of building polygons, heights and lod1 3d models, 2025.
- [10] Nada Tarkhan, Mikita Klimenka, Kelly Fang, Fabio Duarte, Carlo Ratti, and Christoph Reinhart. Mapping facade materials utilizing zero-shot segmentation for applications in urban microclimate research. *Scientific reports*, 15(1):5492, 2025.
- [11] Jinfeng Xie, Minhua Li, Jiaqi Wu, Xiaohu Zhang, and Jie Zhang. Semantic segmentation of building façade materials and colors for urban conservation. *npj Heritage Science*, 13(1), 2025.
- [12] Ying Sun and Zhaolin Gu. Using computer vision to recognize construction material: A trustworthy dataset perspective. *Resources, Conservation and Recycling*, 183:106362, 2022.
- [13] Menglin Dai, Jakub Jurczyk, Hadi Arbabi, Ruichang Mao, Wil Ward, Martin Mayfield, Gang Liu, and Danielle Densley Tingley. Component-level residential building material stock characterization using computer vision techniques. *Environmental science & technology*, 2024.
- [14] Andrey Dimitrov and Mani Golparvar-Fard. Vision-based material recognition for automated monitoring of construction progress and generating building information modeling from unordered site image collections. *Advanced Engineering Informatics*, 28(1):37–49, 2014.
- [15] Jiabin Zhang, Tomohiro Fukuda, and Nobuyoshi Yabuki. Development of a city-scale approach for façade color measurement with building functional classification using deep learning and street view images. *ISPRS International Journal of Geo-Information*, 10(8):551, 2021.

- [16] Teng Zhong, Cheng Ye, Zian Wang, Guoan Tang, Wei Zhang, and Yu Ye. City-scale mapping of urban façade color using street-view imagery. *Remote Sensing*, 13(8):1591, 2021.
- [17] Hao Li, Zhendong Yuan, Gabriel Dax, Gefei Kong, Hongchao Fan, Alexander Zipf, and Martin Werner. Semi-supervised learning from street-view images and openstreetmap for automatic building height estimation.
- [18] Ian Seiferling, Nikhil Naik, Carlo Ratti, and Raphaël Proulx. Green streets – quantifying and mapping urban trees with street-level imagery and computer vision. *Landscape and Urban Planning*, 165:93–101, 2017.
- [19] Fang-Ying Gong, Zhao-Cheng Zeng, Fan Zhang, Xiaojiang Li, Edward Ng, and Leslie K. Norford. Mapping sky, tree, and building view factors of street canyons in a high-density urban environment. *Building and Environment*, 134:155–167, 2018.
- [20] Steffen Knoblauch, Levi Szamek, Jonas Wenk, Iddy Chazua, Innocent Maholi, Maciej Adamiak, Sven Lautenbach, and Alexander Zipf. Uav-assisted municipal solid waste monitoring for informed disposal decisions. In *Proceedings of the 2024 International Conference on Information Technology for Social Good*, pages 105–113, New York, NY, USA, 2024. ACM.
- [21] Jonguk Kim, Hyansu Bae, Hyunwoo Kang, and Suk Gyu Lee. Cnn algorithm for roof detection and material classification in satellite images. *Electronics*, 10(13):1592, 2021.
- [22] Małgorzata B. Starzyńska-Grześ, Robin Roussel, Sam Jacoby, and Ali Asadipour. Computer vision-based analysis of buildings and built environments: A systematic review of current approaches. *ACM Computing Surveys*, 55(13s):1–25, 2023.
- [23] Yecheng Zhang, Huimin Zhao, and Ying Long. Cmab: A multi-attribute building dataset of china. *Scientific data*, 12(1):430, 2025.
- [24] Eike Jens Hoffmann, Yuanyuan Wang, Martin Werner, Jian Kang, and Xiao Xiang Zhu. Model fusion for building type classification from aerial and street view images. *Remote Sensing*, 11(11):1259, 2019.
- [25] OpenMap Development Tanzania. 360° street-level photos: A journey to resilient cities in tanzania using mapillary and openstreetmap., 2024.
- [26] Ali Hatamizadeh, Hongxu Yin, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Global context vision transformers.
- [27] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s.
- [28] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows.

## Appendix



Figure 3: Study area—the Msimbazi River Delta in Dar es Salaam, Tanzania—covered by UAV imagery used throughout the analysis.



Figure 4: Tricycle equipped with a GoPro Max used for collecting 360° SV imagery.

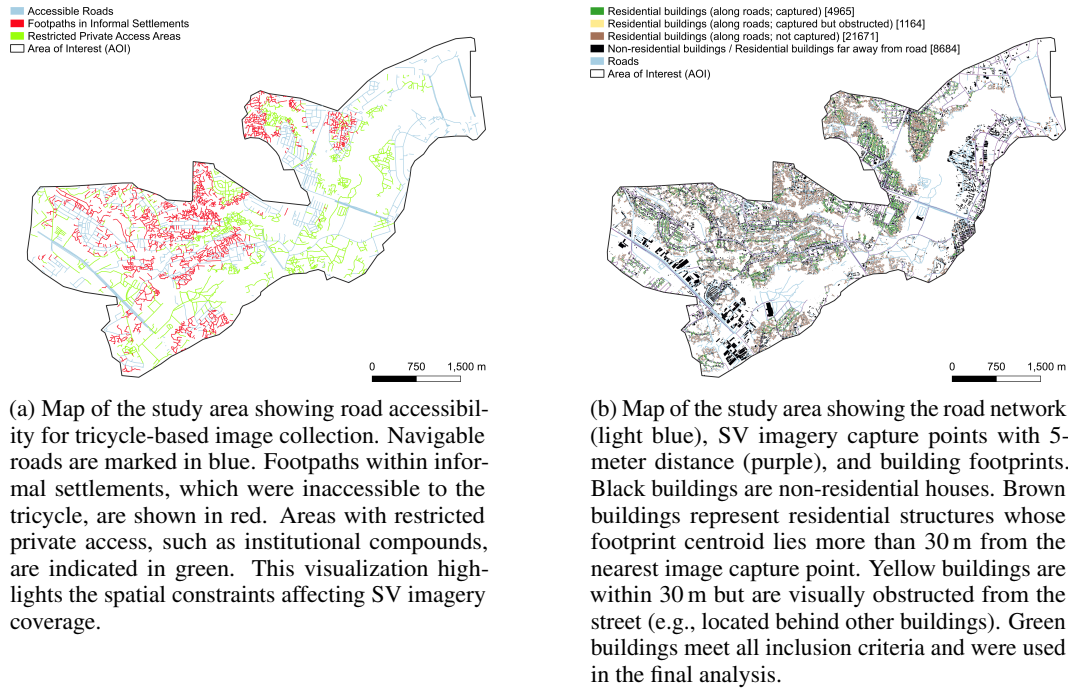


Figure 5: Data acquisition process and spatial coverage of SV imagery. (a) Tricycle-based setup used for image capture. (b) Spatial distribution and visibility classification of buildings in the study area.



Figure 6: Exemplary manually annotated building images showing the range of class labels across different classification tasks.

Table 1: Counts and percentage distribution of annotation classes.

Classification task	Annotation Class	Count	Percentage (%)
Structural openness	Closed Structure	1819	90.95
	Partial	44	2.20
	Unclear	137	6.85
Number of floors	One	1902	95.10
	Two	34	1.70
	Three	11	0.55
	Four+	12	0.60
	Unclear	41	2.05
Vegetation	Yes	1946	97.30
	No	54	2.70
Wall material	Metal	95	4.75
	Concrete	1820	91.00
	Brick	12	0.60
	Wood	5	0.25
	Unclear	68	3.40
Roofing material	Metal	1744	87.20
	Concrete	22	1.10
	Clay	84	4.20
	Tarpaulin	20	1.00
	Wood	63	3.15
	Unclear	67	3.35

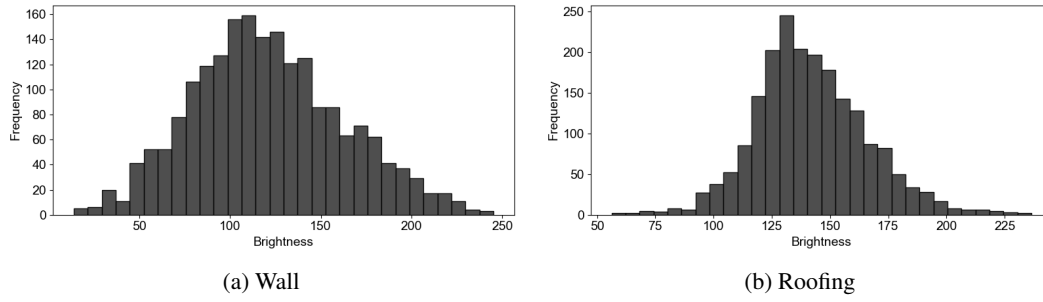


Figure 7: Distribution of brightness values derived from RGB imagery, used as a proxy for surface reflectance and potential heat absorption characteristics.

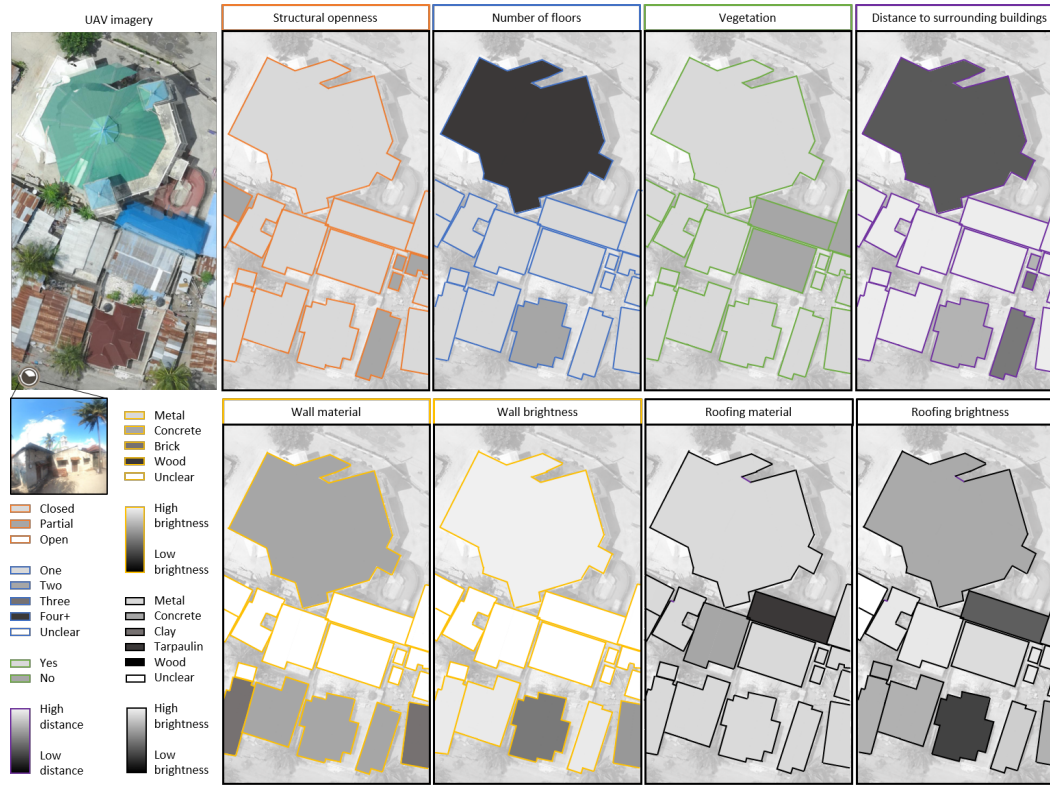


Figure 8: Example of building-level annotations, based on UAV orthophotos and street-view imagery, corresponding to multiple heat-relevant building attributes. Building footprints are outlined and coloured by category: structural openness (orange), number of floors (blue), vegetation (green), distance to surrounding buildings (purple), wall material and brightness (yellow), and roofing material and brightness (black). Each of the eight panels on the right shows a separate attribute label overlaid on the same UAV imagery, enabling multi-label classification of individual buildings for cross-view learning. The imagery was captured near coordinates  $-6.7985, 39.2686$ .

Table 2: Weighted F1-scores (weighted by class support) for classification tasks. Multi-modal results are compared against single-modality models using only UAV or SV data, showing the added value of cross-view learning for building attribute classification. Percentages indicate performance difference relative to multi-modal.

Classification task	Multi-modal F1 (weighted)	SV F1 (weighted)	UAV F1 (weighted)
Vegetation	0.94	0.86 (-9%)	0.80 (-15%)
Number of floors	0.91	0.85 (-7%)	0.45 (-51%)
Roofing material	0.85	0.70 (-18%)	0.82 (-4%)
Wall material	0.68	0.68 (+/-0%)	0.06 (-91%)
Structural openness	0.66	0.57 (-14%)	0.65 (-2%)