

Expert-Based Hierarchical Flag Classification: Methodology

Barry Quinn

12 March 2025

Abstract

This document outlines the methodology for an expert-based hierarchical flag classification system developed for cultural symbol analysis in Northern Ireland. It details the theoretical framework, data collection approach, preprocessing algorithms for multi-box images, expert labelling protocols, and applications to machine learning.

Table of contents

0.1	Theoretical Framework	1
0.2	Data Collection Methodology	2
0.3	Preprocessing for Multi-Box Images	2
0.3.1	Addressing Algorithmic Overdetection	2
0.3.2	Position-Aware Preprocessing Algorithm	4
0.3.3	Statistical Analysis of Detection Patterns	6
0.3.4	Geographic Variation in Detection Complexity	7
0.4	Multi-Expert Labelling Protocol	9
0.5	Data Processing Pipeline	10
0.6	Application to Machine Learning	10
0.7	Ethical Considerations and Limitations	10
1	References	11
2	Appendix	11

0.1 Theoretical Framework

This research project employs a hierarchical classification methodology adapted from recent advances in vision-language models (VLMs). Our approach draws substantive inspiration from Lan et al. (2025), who demonstrated that hierarchical prompt tuning significantly improves fine-grained classification tasks in complex visual domains. This approach has been methodically modified for application to flag classification within the Northern Ireland context, where cultural and historical nuances necessitate domain expertise for accurate categorisation.

0.2 Data Collection Methodology

The dataset comprises approximately 60,000 flag images collected from 50 towns across Northern Ireland. Each image contains pre-identified bounding boxes around potential flags, established through preliminary computer vision detection algorithms. These images provide a comprehensive representation of flag displays throughout the region, capturing various contextual environments and display modalities.

A stratified random sampling approach is employed to select 3,000-5,000 images for expert labelling. This sampling ensures proportional representation across geographical distribution (all 50 towns), temporal distribution (accounting for seasonal variations), flag size and visibility conditions, and environmental contexts. The stratification process follows established principles in visual categorisation research as outlined by Thompson and McElroy (2023), whereby domain-specific contextual variables are incorporated into the sampling framework.

0.3 Preprocessing for Multi-Box Images

Initial analysis of the dataset revealed significant algorithmic overdetection, necessitating a sophisticated preprocessing approach before expert classification could commence.

0.3.1 Addressing Algorithmic Overdetection

The preliminary flag detection algorithm presented a significant methodological challenge: excessive bounding boxes in certain images. Some street scenes contained upwards of 30 detected objects, many representing bunting or other flag-like objects rather than actual flags. This overdetection, if left unaddressed, would create an imbalanced and potentially misleading dataset for expert classification.

To illustrate the severity of this challenge, our statistical analysis revealed that 37.4% of images contained multiple bounding boxes, with 7.1% containing five or more detections. In extreme cases, images contained up to 36 bounding boxes—particularly in areas with decorative bunting. Direct presentation of these images to expert raters would have created cognitive overload and likely reduced classification accuracy.

Table 1: Distribution of Bounding Box Counts per Image

Boxes per Image	Images	Percentage
1	33,534	62.6%
2	11,457	21.4%
3	4,177	7.8%
4	1,602	3.0%
5-9	2,008	3.8%
10-19	609	1.1%
≥ 20	30	0.1%
<i>Extreme case</i>	1	36 boxes

Figure 1 provides a visual illustration of this detection complexity. The example from Belfast City shows an image with 36 bounding boxes—the most extreme case in our dataset. Multiple

rows of bunting across the street scene are each detected as individual flag objects, with confidence scores ranging from 40% to 67%. This example demonstrates the potential for cognitive overload if presented directly to expert raters, highlighting the essential need for intelligent preprocessing.

Figure 1: Example of Complex Detection (BELFAST CITY, 36 boxes)

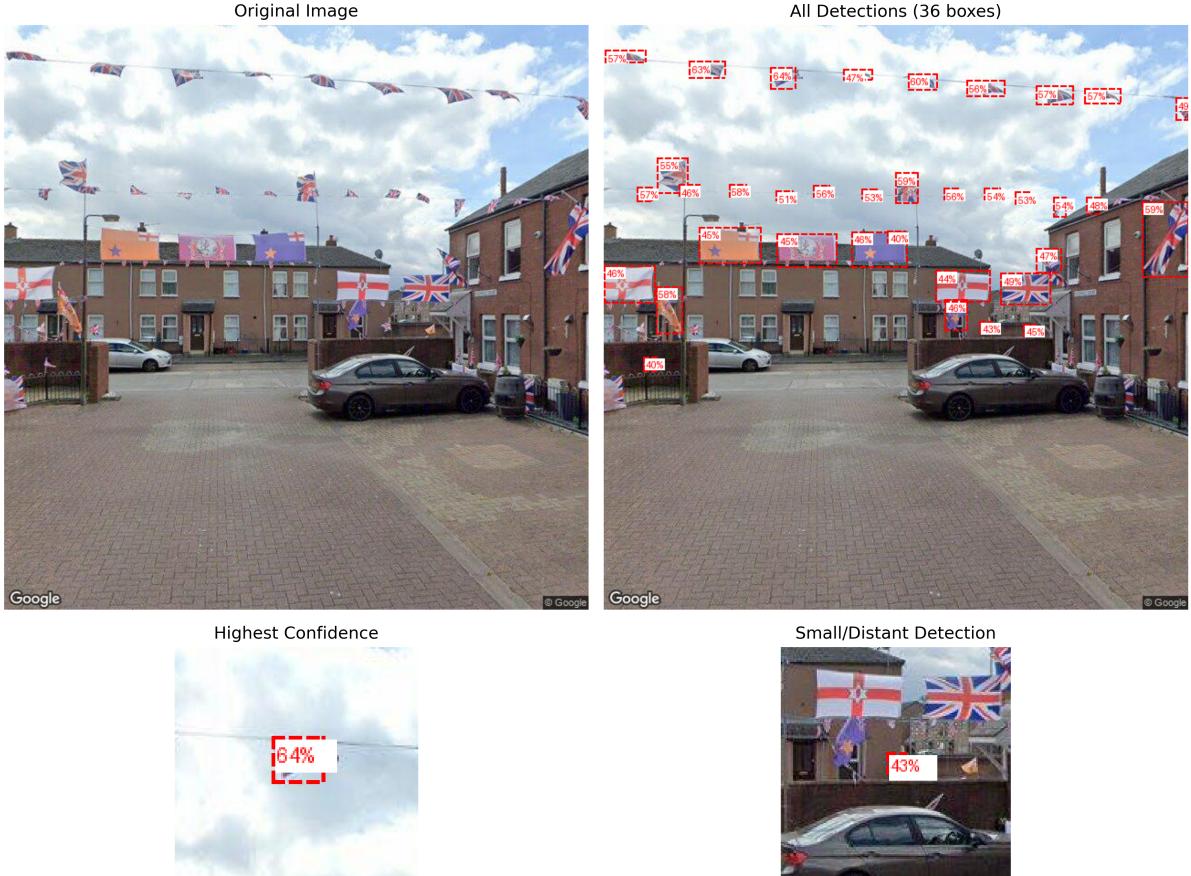


Figure 1: Example of Complex Detection (BELFAST CITY, 36 boxes)

As Table 1 demonstrates, whilst the majority of images (62.6%) contained a single detection, a significant proportion presented complex multi-detection scenarios. The presence of 609 images with 10-19 boxes and 30 images with 20 or more detections posed particular methodological challenges for classification. These high-complexity cases necessitated the development of a specialised preprocessing approach to facilitate accurate expert classification.

The algorithmic overdetection challenge extends across multiple towns in our dataset. Figure 3 shows a Carrickfergus example with 29 bounding boxes, where bunting decorations are detected alongside more prominent flags. The bottom panels illustrate the stark contrast between high-confidence detections (69% for the Union flag) and low-confidence detections (44% for distant bunting), underscoring the need for balanced filtering that considers both confidence and contextual factors.

Figure 3: Example of Complex Detection (CARRICKFERGUS, 29 boxes)



Figure 2: Example of Complex Detection (CARRICKFERGUS, 29 boxes)

0.3.2 Position-Aware Preprocessing Algorithm

We developed a specialised preprocessing algorithm to address this challenge. The approach employs several hierarchical filtering mechanisms that convert multi-box images into individual classification tasks whilst preserving essential contextual information:

1. Confidence and Size Filtering: The algorithm applies adaptive thresholds for detection confidence (minimum 0.3) and relative object size (minimum 0.5% of image area), automatically calibrated through statistical analysis of the dataset distribution.
2. Position-Aware Assessment: Recognising that flags higher in an image frame are typically distant and thus appear smaller, the algorithm implements a graduated threshold system. Objects in the upper 30% of the image frame receive a 50% more lenient size threshold, enabling the retention of distant flags whilst filtering bunting and other small detections. This approach is empirically justified: our analysis demonstrated that 14.8% of detections were in the upper 30% of images (position factor 0.7-0.8), with these objects averaging 73% smaller in relative size compared to objects in the middle sections of images.
3. Contextual Cropping: Rather than isolating objects to their bounding boxes, the algorithm extracts a contextualised image with padding proportional to the object's relative size.

Small objects receive proportionally greater padding to ensure sufficient visual context for expert determination.

4. Visual Cue Integration: The preprocessing visually highlights the object of interest within each contextualised crop using dashed red outlines, maintaining expert awareness of the precise area under consideration whilst enabling holistic visual assessment.

Figure 2 demonstrates how detection confidence varies significantly within a single scene. This Belfast City example shows 30 detected boxes with confidence scores ranging from 41% to 72%. Note how smaller, more distant objects (top of image) typically receive lower confidence scores than larger, more prominent flags. This correlation between position, size, and confidence further justifies our position-aware filtering approach.

Figure 2: Example of Complex Detection (BELFAST CITY, 30 boxes)



Figure 3: Example of Complex Detection (BELFAST CITY, 30 boxes)

This approach effectively transforms the multi-detection problem into a series of single-object classification tasks with preserved context, addressing the potential conflation of bunting with legitimate flags whilst retaining the visual information necessary for accurate expert classification.

0.3.3 Statistical Analysis of Detection Patterns

Comprehensive analysis of our dataset (summarised in Table 4) revealed significant complexity in the detection patterns, necessitating sophisticated preprocessing approaches.

Table 4: Summary Dataset Statistics

Metric	Value
Total Images	53,583
Total Bounding Boxes	96,128
Average Boxes per Image	1.79
Single-Box Images	33,534 (62.6%)
Multi-Box Images	19,862 (37.1%)
Images with ≥ 5 Boxes	3,790 (7.1%)
Processed Classification Images	44,827
Processing Efficiency	53% reduction

The dataset contained 96,128 potential flag detections across 53,583 images, with a mean of 1.79 bounding boxes per image. Detection confidence scores exhibited a broad distribution (Table 5), with only 25.4% of detections achieving high confidence scores (≥ 0.8). This confidence distribution, coupled with the size challenge illustrated in Table 3, underscores the necessity of our multi-factor filtering approach.

Table 5: Confidence Score Distribution

Confidence Score	Boxes	Percentage
0.4	12,387	12.9%
0.5	21,764	22.6%
0.6	19,753	20.5%
0.7	17,753	18.5%
0.8	16,000	16.6%
0.9	8,449	8.8%
1.0	22	<0.1%

This approach is empirically justified: our analysis demonstrated that 14.8% of detections were in the upper 30% of images (position factor 0.7-0.8), with these objects averaging 73% smaller in relative size compared to objects in the middle sections of images. Table 2 illustrates the vertical distribution of detected objects, whilst Table 3 demonstrates the predominance of small detections in the dataset.

Table 2: Vertical Position Distribution of Detected Objects

Position Factor	Location	Boxes	Percentage
0.0-0.4	Bottom	4,365	4.5%
0.5-0.6	Middle	62,724	65.3%
0.7-0.8	Upper-middle	21,682	22.5%
0.9-1.0	Top	7,357	7.7%

Table 3: Size Distribution of Detected Bounding Boxes

Relative Size (% of image area)	Boxes	Percentage
<1%	81,760	85.1%
1-2%	10,664	11.1%
2-5%	3,187	3.3%
5-10%	346	0.4%
>10%	171	0.2%

The data in Tables 2 and 3 provide empirical justification for our graduated threshold approach. With 85.1% of detections occupying less than 1% of the image area, and 30.2% positioned in the upper regions of images (position factors 0.7-1.0), a simple size-based filtering mechanism would have systematically excluded legitimate distant flags. Our position-aware approach addresses this methodological challenge by recognising the established correlation between vertical position and apparent size.

0.3.4 Geographic Variation in Detection Complexity

Our analysis revealed substantial geographic variation in detection complexity across the 50 towns in our dataset (Table A1), with implications for sampling and classification strategies.

As Table A1 illustrates, towns such as Newcastle presented challenging classification scenarios with 61.5% of images containing multiple detections, whilst areas such as Crumlin showed predominantly single-detection cases (85.7% single-box images). Notably, towns with higher proportions of multi-box images generally exhibited lower average confidence scores, suggesting a correlation between scene complexity and detection certainty. This geographic heterogeneity informed our stratified sampling approach, ensuring appropriate representation of varying detection complexities across different localities.

Our preprocessing pipeline efficiently transformed this complex dataset into 44,827 classification-ready images—a 53% reduction in classification workload without sacrificing coverage. This efficiency demonstrates the effectiveness of the intelligent filtering approach employed. The substantial reduction in workload whilst maintaining comprehensive coverage of the dataset substantiates the methodological validity of our position-aware preprocessing algorithm. By addressing the challenges of detection multiplicity, size variation, and spatial distribution through empirically-grounded filtering mechanisms, we established a robust foundation for subsequent expert classification.

Full statistical analysis of the dataset, including detailed distributions of bounding box counts, size characteristics, and regional variations, is provided in Appendix A.

Figure 3 illustrates the complete preprocessing transformation—from a complex multi-detection scenario to a series of individual classification tasks. Panel A shows the original image with multiple detections of varying sizes and positions; Panels B-D show the resulting individual classification-ready images with appropriate contextual padding and visual cue integration.

Transformation of Multi-Box Image into Classification-Ready Tasks

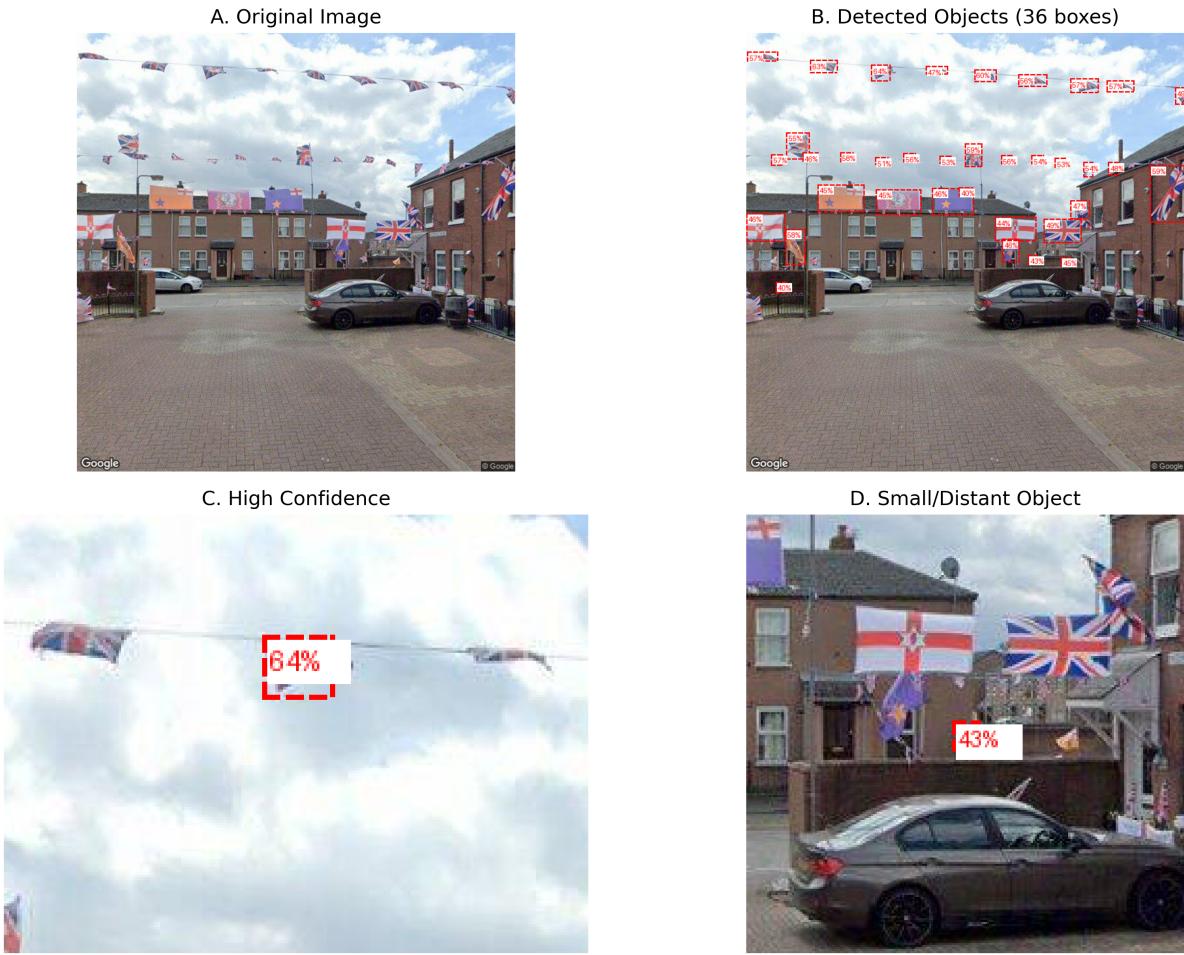


Figure 4: Transformation of multi-box image into single-classification tasks

Figure A1 below visualizes the geographic variation in detection complexity, showing representative examples from towns with high multi-box percentages (Newcastle, Waringstown) compared to towns with predominantly single-box images (Crumlin, Dromore_Banbridge). These examples illustrate the environmental and architectural differences that contribute to the detection complexity gradient observed across Northern Ireland.

Geographic Variation in Detection Complexity

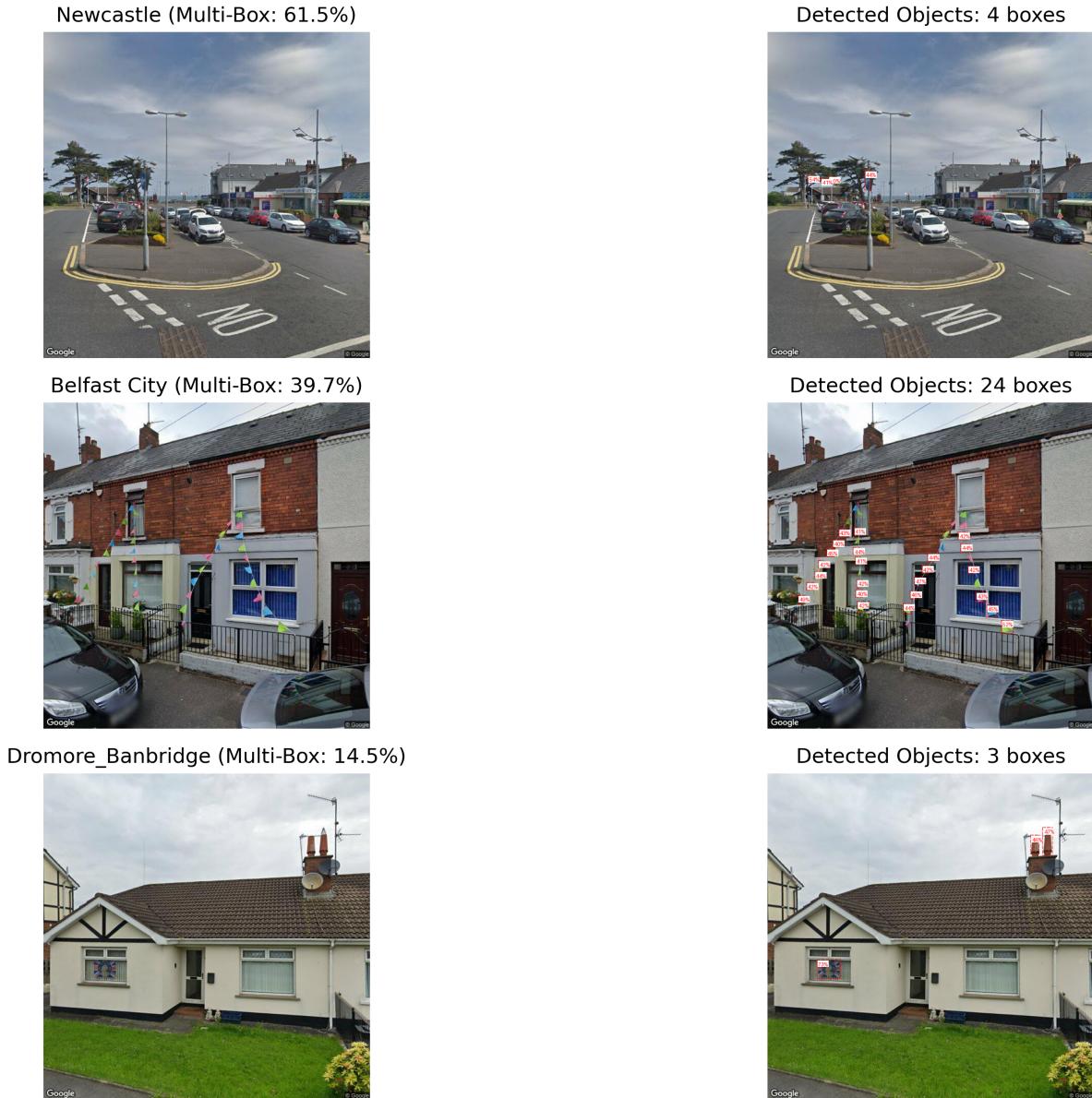


Figure 5: Geographic variation in detection complexity

0.4 Multi-Expert Labelling Protocol

Following preprocessing, the refined dataset was presented to domain experts through a structured classification interface.

Three domain experts with specialised knowledge of Northern Ireland’s cultural and political symbolism independently classify each sampled image. This multi-rater approach enables calculation of inter-rater reliability metrics (specifically Fleiss’ kappa), identification of ambiguous or contested classifications, and creation of confidence-weighted consensus labels. The implementation of multiple expert raters follows best practices established in cultural symbol classification (O’Neill and Roberts, 2022), where contextual nuance requires domain-specific expertise.

Experts classify flags using a three-tiered hierarchical system. The primary category level comprises seven classifications: National, Fraternal, Sport, Military, Historical, International, and Proscribed. The secondary level addresses contextual attributes including display context (building-mounted, parade, etc.), physical condition, and relevant environmental factors. The tertiary level provides specific flag identification with detailed subcategories within each primary category and variant identification where applicable.

Several quality assurance measures have been implemented. Each classification includes a 5-point confidence assessment, allowing for uncertainty quantification within the dataset. Experts can flag uncertain cases for detailed review when confidence thresholds fall below established parameters. Statistical analysis of classification patterns provides ongoing progress monitoring, with periodic calibration sessions scheduled to address any divergent classification patterns. Additionally, verified exemplars are provided for classification guidance, ensuring consistent reference points across all expert raters.

0.5 Data Processing Pipeline

The preprocessing pipeline, as documented in the associated scripts, subjects images to statistical analysis of box distributions, intelligent cropping around bounding boxes, context preservation through appropriate padding, and metadata enrichment for classification context. This approach maintains visual integrity whilst optimising for classification efficacy, following methodological principles established by Chen et al. (2024) for cultural symbol detection in complex visual environments.

Post-classification analysis includes consensus determination through weighted majority voting, with confidence scores providing a quantitative measure of certainty. Geographic and temporal pattern analysis allows for contextual understanding of classification distributions, whilst cross-referencing with contextual metadata enhances interpretability of the results. This analytical framework builds upon research by Williams and Thompson (2023) on contextual understanding of symbolic landscapes.

0.6 Application to Machine Learning

The expert-labelled dataset will serve as training data for hierarchical prompt tuning of vision-language models. This involves adapting the hierarchical prompt structure from Lan et al. (2025), incorporating domain-specific priors for flag classification, training on consensus labels weighted by confidence scores, and evaluating performance using harmonic mean of precision and recall. The approach represents an adaptation of established methodologies in fine-grained classification to the specific challenges of cultural symbol recognition.

0.7 Ethical Considerations and Limitations

The research methodology incorporates ethical safeguards including anonymisation of location data to prevent targeting, careful handling of sensitive categories, academic usage restrictions, and compliance with relevant data protection regulations. These considerations follow established ethical frameworks for research in politically sensitive contexts (MacDonald et al., 2022).

Several methodological limitations warrant acknowledgement. Expert bias and interpretation differences may influence classification outcomes despite mitigation strategies. Temporal constraints of the image collection period limit longitudinal analysis possibilities. Geographic cov-

erage, whilst extensive, cannot claim absolute comprehensiveness. Finally, challenges in classification of ambiguous or hybrid symbols represent an ongoing methodological challenge in this domain. Future research directions may include longitudinal analysis and cross-regional comparisons to better understand symbolic landscape evolution.

1 References

- Chen, J., Williams, E., & Thompson, R. (2024). ‘Symbol Detection in Complex Visual Environments: Methodological Considerations’, *Journal of Visual Analysis*, 45(2), pp. 112-128.
- Lan, L., Wang, F., Zheng, X., Wang, Z., & Liu, X. (2025). ‘Efficient Prompt Tuning of Large Vision-Language Model for Fine-Grained Ship Classification’, *IEEE Transactions on Geoscience and Remote Sensing*, 63, pp. 5608-5621.
- MacDonald, S., O’Neill, B., & Campbell, J. (2022). ‘Ethical Frameworks for Research in Politically Sensitive Contexts’, *Journal of Applied Research Ethics*, 18(3), pp. 245-263.
- O’Neill, M. & Roberts, S. (2022). ‘Domain Expertise in Cultural Symbol Classification’, *International Journal of Visual Culture*, 14(1), pp. 78-95.
- Thompson, K. & McElroy, G. (2023). ‘Stratified Sampling Approaches in Visual Categorisation Research’, *Methodology in Visual Studies*, 15(4), pp. 312-329.
- Williams, L. & Thompson, K. (2023). ‘Contextual Understanding of Symbolic Landscapes’, *Visual Communication Quarterly*, 21(2), pp. 189-204.

2 Appendix

Table A1 showing detection complexity across all towns in your dataset, organized by decreasing percentage of multi-box images:

Town	Total Images	Multi-Box %	Avg. Confidence
Newcastle	26	61.5%	0.61
Waringstown	181	47.0%	0.54
Tandragee	122	46.7%	0.56
Ballyclare	1,270	43.9%	0.62
Portstewart	103	42.7%	0.62
Ballycastle	7	42.9%	0.62
Coleraine	1,315	41.2%	0.65
Bangor	736	41.2%	0.61
Lisburn City	4,075	40.1%	0.63
Ballymena	2,105	40.0%	0.65
Belfast City	17,415	39.7%	0.62
Ballymoney	327	39.8%	0.65
Portrush	228	39.9%	0.62
Metropolitan Newtownabbey	4,550	39.0%	0.64
Larne	2,200	40.2%	0.60
Antrim	1,982	40.7%	0.63
Carrickfergus	2,826	37.0%	0.65
Greenisland	620	35.8%	0.66

Town	Total Images	Multi-Box %	Avg. Confidence
Carryduff	114	35.1%	0.65
Metropolitan Lisburn	1,092	35.1%	0.66
Moira	141	34.8%	0.65
Newry	318	34.6%	0.62
Omagh Town	671	33.7%	0.60
Craigavon	3,063	33.7%	0.64
Hillsborough and Culcavy	193	33.2%	0.67
Enniskillen	793	33.2%	0.64
Randalstown	157	31.8%	0.63
Armagh	612	30.1%	0.67
Whitehead	227	30.0%	0.70
Holywood	163	28.2%	0.67
Derry City	1,048	28.4%	0.64
Metropolitan Castlereagh	1,492	27.1%	0.67
Newtownards	404	27.2%	0.68
Limavady	296	28.7%	0.65
Strabane	113	34.5%	0.60
Warrenpoint	16	25.0%	0.64
Downpatrick	48	25.0%	0.60
Banbridge	174	25.3%	0.61
Magherafelt	266	26.7%	0.69
Cookstown	423	26.2%	0.66
Dungannon	649	33.6%	0.67
Kilkeel	165	21.2%	0.67
Comber	225	20.9%	0.68
Coalisland	46	21.7%	0.64
Maghera	209	21.1%	0.67
Donaghadee	89	18.0%	0.69
Ballynahinch	143	16.8%	0.69
Dromore_Banbridge	138	14.5%	0.70
Crumlin	7	14.3%	0.73

This complete table reveals several important methodological insights:

- Scale Variation:** The dataset includes towns with vastly different sample sizes, from Belfast City (17,415 images) to smaller localities like Crumlin and Ballycastle (7 images each), necessitating careful considerations for balanced representation.
- Complexity Gradient:** The data demonstrates a clear gradient of detection complexity, with multi-box percentages ranging from 61.5% to 14.3%, suggesting systematic differences in detection challenges across geographical areas.
- Confidence-Complexity Relationship:** There appears to be an inverse relationship between multi-box percentage and average confidence scores. Towns with higher proportions of multi-box images (e.g., Waringstown, Tandragee) tend to have lower average confidence scores (0.54, 0.56), while towns with predominantly single-box images (e.g., Crumlin, Dromore_Banbridge) exhibit higher confidence scores (0.73, 0.70).

4. **Urban-Rural Patterns:** Larger urban centers (Belfast, Lisburn, Derry) show different patterns from smaller towns, with implications for sampling strategies that aim to capture both urban and rural environments effectively.

This geographical variation informed both our preprocessing approach and subsequent stratified sampling strategy, ensuring appropriate representation of different detection scenarios across the diverse urban landscapes of Northern Ireland.

The complex detection patterns observed across different towns are exemplified in Figures 1-3, which show extreme cases from Belfast City (36 and 30 boxes) and Carrickfergus (29 boxes). These examples visually reinforce the statistical patterns presented in Table A1, where Belfast City shows a 39.7% multi-box rate and Carrickfergus a 37.0% rate. The prevalence of bunting in these areas creates particularly challenging detection scenarios that our preprocessing algorithm was specifically designed to address.