
BETA DISTRIBUTION LEARNING FOR RELIABLE ROADWAY CRASH RISK ASSESSMENT

Ahmad Elallaf¹ Nathan Jacobs² Xinyue Ye³ Mei Chen⁴ Gongbo Liang¹

¹Texas A&M University-San Antonio ²Washington University in St. Louis

³University of Alabama ⁴University of Kentucky

{aelallaf, gliang}@tamus.edu jacobsn@wustl.edu xye10@ua.edu mei.chen@uky.edu
<https://www.gb-liang.com/projects/betarisk>

ABSTRACT

Roadway traffic accidents represent a global health crisis, responsible for over a million deaths annually and costing many countries up to 3% of their GDP. Traditional traffic safety studies often examine risk factors in isolation, overlooking the spatial complexity and contextual interactions inherent in the built environment. Furthermore, conventional Neural Network-based risk estimators typically generate point estimates without conveying model uncertainty, limiting their utility in critical decision-making. To address these shortcomings, we introduce a novel geospatial deep learning framework that leverages satellite imagery as a comprehensive spatial input. This approach enables the model to capture the nuanced spatial patterns and embedded environmental risk factors that contribute to fatal crash risks. Rather than producing a single deterministic output, our model estimates a full Beta probability distribution over fatal crash risk, yielding accurate and uncertainty-aware predictions—a critical feature for trustworthy AI in safety-critical applications. Our model outperforms baselines by achieving a 17-23% improvement in recall, a key metric for flagging potential dangers, while delivering superior calibration. By providing reliable and interpretable risk assessments from satellite imagery alone, our method enables safer autonomous navigation and offers a highly scalable tool for urban planners and policymakers to enhance roadway safety equitably and cost-effectively.

1 Introduction

Roadway traffic accidents claim over 1.3 million lives annually [1] and impose economic burdens of 3% of the GDP in many countries [2]. As a critical infrastructure sector [3], transportation safety has garnered significant research [4–6], yet accurately estimating crash risk remains a challenge due to its inherent uncertainties and the sparse nature of crash events.

Conventional safety research often analyzes individual factors separately, such as driver behavior [7], road infrastructure [8], traffic patterns [9], and weather [10], overlooking the complex interplay between these elements [11]. Since crash occurrences frequently result from intricate multi-factor interactions, methods that analyze these factors in isolation struggle to predict risk holistically [12]. Furthermore, data limitations have constrained the scope of most studies to highways [13–17], leaving comprehensive crash risk analysis for local roads, where data is often less available, relatively unexplored.

To overcome these limitations, we introduce a novel deep learning framework that learns a full Beta probability distribution, moving beyond simple point-estimates of fatal crash risk. Our primary contributions are threefold:

- A **holistic, vision-based model** that captures the complex interplay of risk factors embedded in the visual data, in contrast to methods that study variables in isolation.
- A **probabilistic formulation** that yields well-calibrated, uncertainty-aware predictions, a critical feature for trustworthy AI in high-stakes, safety-critical domains.

- A **highly scalable and equitable methodology** that uses near-globally available satellite imagery, enabling risk assessment for both highways and previously under-assessed local roads.

The proposed probabilistic model is evaluated through extensive experiments conducted over four major metropolitan areas, which have a population of ≈ 20 million. Our model achieves a 17–23% improvement in recall over baselines, a crucial metric for any safety-critical task, while also delivering superior model calibration and F1 scores. By producing reliable and interpretable risk assessments from satellite imagery alone, this work provides a foundational tool for enhancing traffic safety, from enabling safer route selection for drivers and autonomous vehicles to empowering urban planners and policymakers to mitigate high-risk areas.

2 Background

2.1 Estimate Roadway Crash Risk

A primary challenge in data-driven roadway safety is formulating the risk estimation task. Existing methods often frame it as classifications, such as predicting a crash occurrence within a short time frame [9]. While valuable, these approaches do not estimate the inherent, continuous crash risk of a given road segment. A more nuanced approach is to directly estimate a crash probability, such as using Monte Carlo simulations [18–20]. However, this is fundamentally challenged by the extreme sparsity of crash data. For instance, the average annual accident rate for a 25m^2 road segment in the United States is just 0.1% [21]. This level of sparsity renders traditional estimation techniques unreliable, as they can obscure high-risk areas while falsely flagging safe ones [22], leading to false negatives that are dangerous in any safety-critical application. Furthermore, such simulation methods are often ill-suited for large-scale applications due to high computational costs and the need for carefully tuned parameters.

Deep Neural Networks (DNNs) offer a powerful alternative, as they can learn complex, task-specific features directly from data and provide near-instantaneous inference. However, supervised DNNs typically rely on large, manually labeled datasets, such as manually assigned risk levels (e.g., low, neutral, high) [23]. Creating these datasets is prohibitively expensive, and the manual labels can suffer from human bias, potentially misrepresenting the true risk [24, 25]. These challenges motivate the need for a new approach that can learn a continuous risk score from objective crash data while effectively handling the probabilistic nature of the task.

2.2 Deep Neural Network Miscalibration

Over the recent years, DNNs have shown promising performance on various domains, such as medical imaging [26, 27], cybersecurity [28], transportation [29], and astrophysics [30]. However, for a predictive model to be trustworthy in high-stakes applications, its predicted confidence must accurately reflect its probability of being correct. However, modern DNNs are often miscalibrated, tending to produce overconfident predictions [31, 32].

Mathematically, a model is perfectly calibrated if, for any given confidence level p , the long-run accuracy of predictions with that confidence is indeed p . For DNNs, the calibration error, the difference between a model’s predicted confidence and its actual accuracy, is often significantly greater than zero [33]. This miscalibration is a critical failure point in high-stakes applications where decisions depend on the model’s self-assessed certainty.

While various techniques can mitigate this issue, they often have limitations. Post-processing methods like temperature scaling [32] adjust model outputs without altering the learned features, while in-training regularization [34, 35] requires careful tuning for the weight scaler. Given that model complexity is a key contributor to miscalibration [36], we argue that an effective solution must be deeply integrated into the learning process. Our work achieves this by reformulating the risk estimation task as learning a full probability distribution, a method that inherently encourages better-calibrated and more reliable predictions.

3 Method

3.1 Probabilistic Modeling Framework

Our method recasts roadway crash risk estimation from a standard classification task into a probabilistic learning problem, motivated by the limitations of conventional models that provide a single point-estimation. Consider a fatal crash, a stochastic occurrence, at a specific point in spacetime, $C = (x, y, t, d)$, where (x, y) is the geolocation and (t, d) is the time and date. While any single crash is a random event, its location provides the strongest available evidence for a local maximum in the underlying, continuous risk field, $R(\cdot)$. Therefore, it is intuitive that the inferred risk should be higher at or near the crash site and should decay smoothly as one moves away in space or time. For

Algorithm 1 Target Beta Distribution Generation

Require: Original image x , binary label $l \in \{0, 1\}$, base concentration K_{base} , minimum positive risk mean μ_{min} , minimum positive concentration k_{min} , distance weight w_{dist} , size weight w_{size} , and $\epsilon = 1e^{-5}$

```

if  $l = 0$  then                                 $\triangleright$  For negative samples, create a low-risk,
     $\alpha_t \leftarrow \epsilon$                           $\triangleright$  high-certainty distribution
     $\beta_b \leftarrow K_{base}$ 
else                                          $\triangleright$  For positive samples ( $l=1$ ), generate labels
     $x' \leftarrow$  random crop of  $x$             $\triangleright$  based on crop geometry
     $d_{norm} \leftarrow$  normalized distance of  $x'$  from center of  $x$ 
     $s_{norm} \leftarrow \frac{\text{size}(x')}{\text{size}(x)}$ 
     $influence \leftarrow w_{dist} \cdot (1 - d_{norm}) + w_{size} \cdot s_{norm}$ 
     $\mu_t \leftarrow \mu_{min} + (1 - \mu_{min}) \cdot influence$ 
     $k_t \leftarrow k_{min} + (K_{base} - k_{min}) \cdot influence$ 
     $\alpha_t \leftarrow \mu_t \cdot k_t$ 
     $\beta_t \leftarrow \epsilon$ 
end if
return  $(\alpha_t, \beta_t)$                                  $\triangleright$  Return the target Beta distribution

```

nearby points, such as a spatially displaced point $C' = (x - \delta, y, t, d)$, the risk should be lower, i.e., $R(C') < R(C)$. Standard point-estimate classifiers fail to capture this continuous field, as they are trained to predict a binary outcome for each location independently.

While a complete model would account for both spatial and temporal decay, this work focuses on the challenging and foundational task of estimating the **static, inherent risk** of a location based on its geographic and structural features. Our goal is to model the spatial component of this uncertainty by learning a distribution over possible risk values, capturing the intuition that for a nearby point C' , the risk is attenuated but non-zero: $0 < R(C') < R(C)$.

We specifically employ the Beta distribution for this task due to its natural support on the $[0, 1]$ interval and its flexibility in representing diverse risk profiles. Instead of a single value, our model $h(x)$ maps an input image x to the two positive scalar parameters, (α, β) , which define a Beta distribution, $P_p \sim Beta(\alpha, \beta)$. This formulation allows the model to express its uncertainty through the shape of the distribution: a sharp peak indicates high confidence, while a wide distribution signifies high uncertainty. The final risk score R is the mean of this predicted distribution:

$$R = \mathbb{E}[P_p] = \frac{\alpha}{\alpha + \beta}. \quad (1)$$

To achieve this, our framework integrates three key technical contributions: 1) a novel procedural labeling technique that generates the targeting Beta distributions from data augmentation, 2) a multi-scale deep neural network architecture, and 3) a compound loss function for joint optimization.

3.2 Target Beta Distributions Generation

A key innovation of our framework is the procedural generation of supervisory signals in the form of target Beta distributions. Instead of using static labels, we dynamically create a target Beta distribution, $P_t \sim Beta(\alpha_t, \beta_t)$, for each training sample based on the properties of the random crop augmentation. Specifically, given an input image, we first apply a random crop. The target Beta distribution is, then, generated using Algorithm 1. This process acts as a sophisticated form of structured label smoothing, transforming data augmentation from a simple regularizer into a rich source of continuous supervision for risk and uncertainty.

For **negative samples** (no crash), the objective is to predict low risk with high confidence. The target distribution is therefore constant: α_t is set to a small positive value ϵ and β_t is set to a large value representing high certainty K_{base} , creating a distribution sharply peaked at zero.

For **positive samples** (crash), the target distribution reflects the quality of the visual evidence in the random crop. This is quantified by an **influence** score, which modulates the target distribution's mean and concentration to generate a supervision signal that is proportional to the information content of the augmented image. The score is a weighted combination of two geometric properties of the crop: its centrality relative to the crash location and its size.

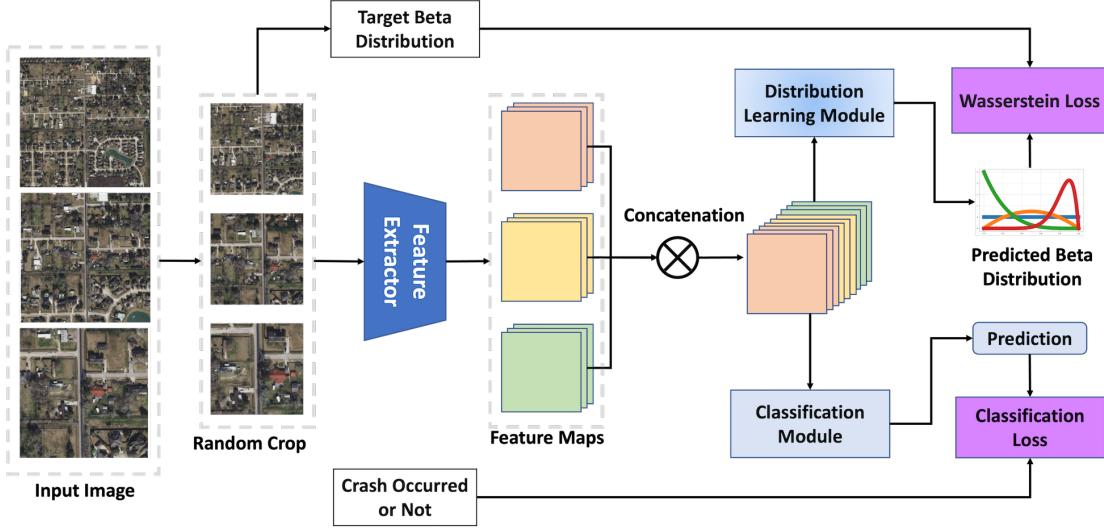


Figure 1: Training Architecture with Joint Optimization

We set the weights to 0.7 for centrality (w_{dist}) and 0.3 for relative size (w_{size}). This weighting scheme is based on the strong intuition that the visual features most critical to understanding risk—such as specific road geometry, lane markings, the presence of an intersection, or the surrounding environments—are spatially concentrated around the event’s location. A crop that is well-centered on the crash point provides the clearest and most relevant evidence, thus deserving a higher influence score and a more confident target distribution. The relative size of the crop provides useful, but secondary, broader context about the surrounding environment. This principled approach transforms data augmentation into a rich source of supervision, teaching the model to dynamically associate higher risk and confidence with visual samples that contain the most informative evidence.

This influence score then modulates the target mean μ_{min} and the target concentration k_t , which in turn define the final Beta parameters. For positive samples, the β_t is set to the small constant ϵ , ensuring the distribution is always skewed towards high risk, with the influence score controlling the precise shape and confidence (see the **supplement materials** for the hyperparameters used in this study).

3.3 Model Architecture

The architecture of our model, illustrated in Figure 1, is designed to process multi-scale satellite imagery. During training, a random crop is sampled from the input, which consists of image slices of the same location at different resolutions.

The cropped images are, then, passed through a shared feature extractor backbone to produce multiple corresponding feature maps. These maps are concatenated along the channel dimension to form a unified feature representation, serving as the input for two parallel prediction heads:

- A Distribution Learning Head, which outputs the two Beta parameters (α, β) .
- An auxiliary Classification Head, which outputs a single logit for the binary crash/no-crash task.

3.4 Training and Optimization

The model is trained end-to-end by jointly optimizing the two parallel heads with a compound loss function. The primary distribution learning head is supervised by the a mean-variance loss that inspired by the squared Wasserstein-2 (W_2^2) distance [37], which measures the dissimilarity between the predicted (P_p) and the target (P_t) Beta distributions:

$$\mathcal{L}_{W_2^2}(P_p, P_t) = (\mu_p - \mu_t)^2 + (\sigma_p - \sigma_t)^2, \quad (2)$$

where the μ and σ are the mean and standard deviation.

We empirically selected this W_2^2 surrogate over true W_2^2 distance and other distribution divergence metrics, including KL-Divergence [38] and the Cramér-von Mises criterion [39]. As a true metric, our W_2^2 surrogate loss provides a

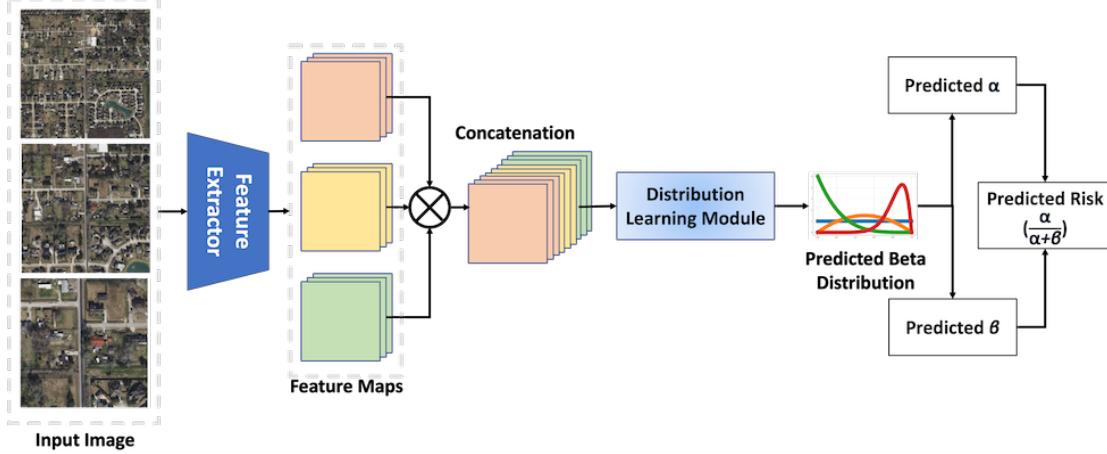


Figure 2: Streamlined Architecture for Inference

more stable gradient than KL-Divergence, especially when the predicted and target distributions have little overlap. Most importantly, for one-dimensional distributions like the Beta, the W_2^2 surrogate loss directly optimizes of the risk score (the mean) and confidence level (the standard deviation) simultaneously. Our experimental analysis also shows this surrogate is a close approximation of the true W_2^2 (errors on the order of 10^{-3} to 10^{-2}), deviating only in extreme cases (see the **supplementary materials**).

The auxiliary classification head is supervised by a Binary Cross-Entropy loss, which encourages the shared backbone to learn discriminative features relevant to the safety task:

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)], \quad (3)$$

where y_i and p_i are the label and predicted probability.

The overall objective function is a weighted combination of the two losses, balanced by hyperparameters, λ_1 and λ_2 :

$$\mathcal{L} = \lambda_1 \cdot \mathcal{L}_{BCE} + \lambda_2 \cdot \mathcal{L}_{W_2^2}. \quad (4)$$

3.5 Inference Process

The inference process, illustrated in Figure 2, is direct and computationally efficient. The random crop augmentation and the auxiliary classification head are removed. The full, uncropped multi-scale image is passed through the feature extractor backbone and the distribution head. The risk score R is calculated as the mean of the distribution, per Equation 1. This feed-forward process allows for rapid and scalable risk assessment of any location.

4 Experiment Setup

This study utilizes the MSCM dataset [40], a large-scale collection of multi-scale satellite images from Texas, USA, with 16,451 locations labeled with historical fatal crashes. All models use a ResNet-50 [41] backbone, $\lambda_1 = 5$, $\lambda_2 = 1$, and were trained on NVIDIA A100 GPUs. See the **supplementary materials** for more information about the dataset, implementation details, and hyperparameter analysis for the selection of λ_1 and λ_2 .

4.1 Evaluation Methodology

Quantitative Metrics We first evaluate our model’s practical effectiveness by framing the risk estimation as a binary classification task to identify historical crash locations. The model’s predicted risk score R , derived from Equation 1 is thresholded at 0.5 to yield a binary prediction. We then assess the model’s predictive performance using standard metrics: F1-Score, Precision, Recall, AUC (Area Under the Receiver Operating Characteristic curve), and PRC (area under the precision-recall curve); and assess model’s calibration using Expected Calibration Error (ECE) and Brier score. Due to safety-oriented, we consider **Recall to be the most critical metric** that answers the question: “*Of all crash locations, what fraction did our model successfully identify?*”

Table 1: Main Quantitative Results (**bold**: best performance; underlined: second best performance)

Methods	Pre-Train	Probabilistic	Multi-Scale	Performance (\uparrow)					Uncertainty (\downarrow)	
				F1	Precision	Recall	AUC	PRC	ECE	Brier
ImageNet	ImageNet	x	x	0.4753	0.4968	0.4555	0.7980	0.4862	0.1281	0.1600
MSCM-SS	MSCM	x	x	0.4966	0.4981	0.4950	0.8165	0.5185	<u>0.1006</u>	0.1458
MSCM-MS	MSCM	x	✓	0.5409	0.6731	0.4521	0.8572	0.6269	0.1067	0.1296
Prob-SS (Ours)	MSCM	✓	x	0.5001	0.4252	0.6070	0.7749	0.4409	0.1731	0.1922
Prob-MS (Ours)	MSCM	✓	✓	0.5762	0.6296	<u>0.5311</u>	0.8663	0.6489	0.0881	0.1211

We also evaluate our method against a Deep Ensemble (DE) of the strongest baseline, constructed from three independent training runs. The final predicted risk score of a DE model, R_{DE} , is calculated as the mean of the predictions from each individual model in the ensemble. This single score for each sample is then used to compute all the aforementioned performance and calibration metrics.

The ensemble’s predictive uncertainty is quantified in two ways: the variance of the risk scores and the disagreement rate among the final binary predictions. A higher value in either metric reflects greater disagreement among the models and thus higher uncertainty in the final prediction.

Qualitative Analysis To intuitively understand the value of our probabilistic approach, we conduct a qualitative analysis of the model’s outputs from two perspectives. First, we analyze the aggregate behavior of the model’s outputs by comparing the overall distribution of predicted probabilities from our model against the baselines. By plotting a histogram of all risk scores, we can visually assess model confidence. A well-calibrated, uncertainty-aware model is expected to utilize the full [0, 1] probability range, whereas overconfident models will show predictions heavily clustered at the extremes (near 0 and 1). Second, we visualize the predicted Beta distributions for four distinct scenarios: true positives (TP), true negatives (TN), false negatives (FN), and false positives (FP). The goal of this analysis is to provide an intuitive understanding of the model’s behavior by interpreting its successes and failures.

Case Study: San Antonio River Walk To demonstrate the model’s utility, we conduct a case study of the San Antonio River Walk, providing practical insight into the model’s performance in a challenging, safety-critical area.

4.2 Baseline Models

We evaluate our method against three baselines to isolate our framework’s contributions. Our primary benchmark is the Multi-Scale Cross-Matching (MSCM) model [40], the current state-of-the-art for fatal crash risk estimation using only satellite imagery.

ImageNet Baseline: A standard model pre-trained on ImageNet [42] that takes single-scale satellite images as input, providing us the performance of a generic, non-domain-specific feature extractor on our task.

MSCM-SS (Single-Scale): The same single-scale architecture but using weights generated by the self-supervised pre-training through cross-matching, proposed in the MSCM paper, to test the value of domain-specific features.

MSCM-MS (Multi-Scale): The full MSCM model, which uses both its domain-specific pre-training and multi-scale imagery as input, represents the strongest available baseline, allowing us to compare against the current state-of-the-art classification approach directly.

The best checkpoint for each model was selected based on the model accuracy on the validation set.

5 Results

5.1 Quantitative Analysis

Table 1 summarizes the quantitative results. The single-scale baselines (ImageNet and MSCM-SS) achieve a < 0.5 precision and recall scores, indicating their predictions for positive cases are close to random and exhibit little ability to identify high-risk areas. While the MSCM-MS model achieves high precision (0.6731), its poor recall (0.4521) means it fails to identify over half of all crash locations, rendering it unreliable for safety-critical applications.

Table 2: Deep Ensemble Results over Three Training Trails (**bold**: best performance)

Methods	Performance (\uparrow)					Uncertainty (\downarrow)		Disagreement (\downarrow)	
	F1	Precision	Recall	AUC	PRC	ECE	Brier	Variance	Disagr. Rate
Ensemble MSCM-MS	0.5966	0.7062	0.5165	0.8839	0.6890	0.0787	0.1112	0.0925	16.93%
Ensemble Prob-MS (Ours)	0.5976	0.6750	0.5361	0.8761	0.6886	0.0605	0.1075	0.0822	15.14%

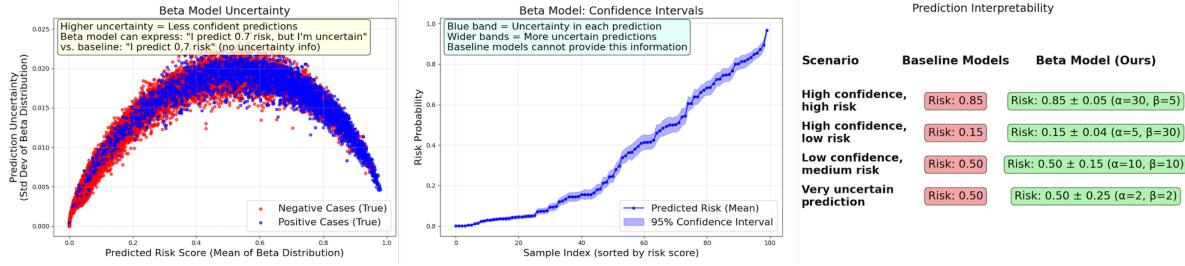


Figure 3: Uncertainty Quantification and Interpretability

In contrast, our models demonstrate a significant improvement in identifying potential dangers. Our multi-scale model, Prob-MS, achieves the best overall balance of performance, attaining the highest F1-score. Its most significant contribution is boosting the recall to 0.5311, a 17% relative improvement over MSCM-MS, drastically reducing the number of hazardous sites that would be missed. Our single-scale model, Prob-SS (0.6070 recall score), significantly improves the metric by 23% over the best baseline (0.4950).

Crucially, Prob-MS is also the most trustworthy model, achieving the lowest (best) ECE of 0.881 and Brier of 0.1211. This confirms that our model's probabilistic outputs are more statistically sound and reliable for real-world decision-making.

We also evaluate our method against a Deep Ensemble of the strongest baseline (Table 2). When comparing our single Prob-MS model against the baseline Ensemble MSCM-MS, we find that our single model achieves competitive performance, including a 3% higher recall, better calibration, and lower uncertainty at only 1/3 the computational cost at both training and inference times. This highlights the efficiency and practical advantage of our approach.

In an apples-to-apples comparison between ensembled methods, our Ensemble Prob-MS demonstrates the clear superiority of our probabilistic framework. It outperforms the baseline ensemble on the most critical metrics for this task, achieving a higher F1-score and recall. Most importantly, it is significantly better calibrated and exhibits lower uncertainty, as evidenced by its superior (lower) ECE, Brier, Variance, and Disagreement Rate scores.

5.2 Qualitative Analysis

Our qualitative analysis highlights the superior interpretability and trustworthiness of our probabilistic framework. As shown in Figure 3, our model provides a comprehensive and practical understanding of risk that standard classifiers cannot offer. The “Beta Model Uncertainty” plot (left) confirms the model’s rational behavior, showing that prediction uncertainty is lowest for highly confident predictions and highest for ambiguous ones around a 0.5 risk score. The “Confidence Intervals” plot (center) demonstrates that every prediction is accompanied by a 95% confidence interval, with the interval’s width directly communicating the model’s certainty on a per-prediction basis. Finally, the “Prediction Interpretability” table (right) crystallizes this key advantage, showing how our Beta model resolves the ambiguity of a baseline’s “Risk: 0.50” output by distinguishing between a low-confidence prediction (e.g., with $\alpha = 10, \beta = 10$) and a very uncertain one (e.g., with $\alpha = 2, \beta = 2$). This additional context is invaluable for any safety-critical application.

This nuanced, per-prediction behavior leads to a more rational distribution of predictions in aggregate (Figure 4). While baseline models behave like overconfident black boxes with predictions heavily clustered at the extremes of 0 and 1, our model utilizes the full probability spectrum to express varying degrees of certainty. This ability to be “less confident” is not a weakness but a hallmark of a more honest and trustworthy risk assessment tool.

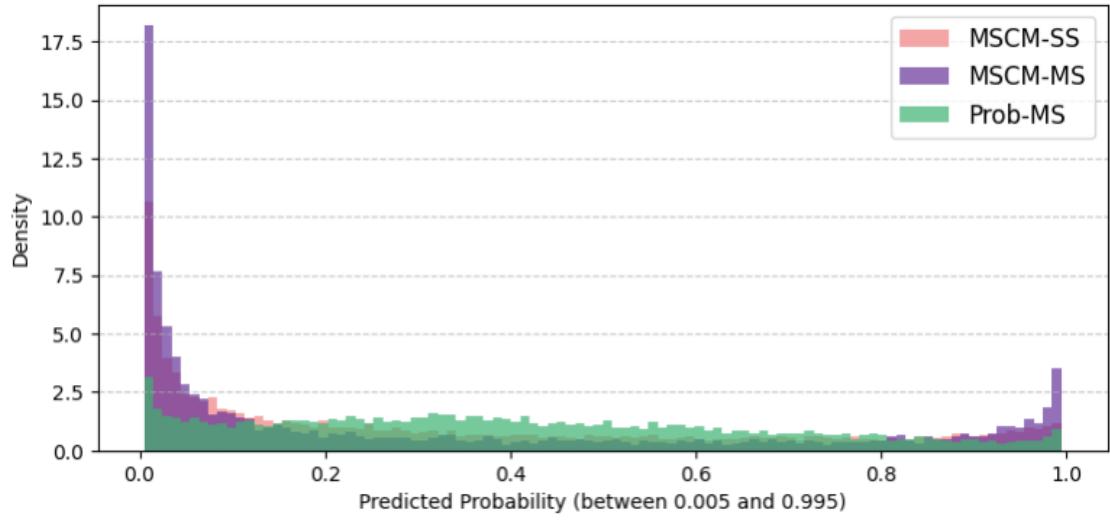


Figure 4: Analysis of Predicted Probability Distributions

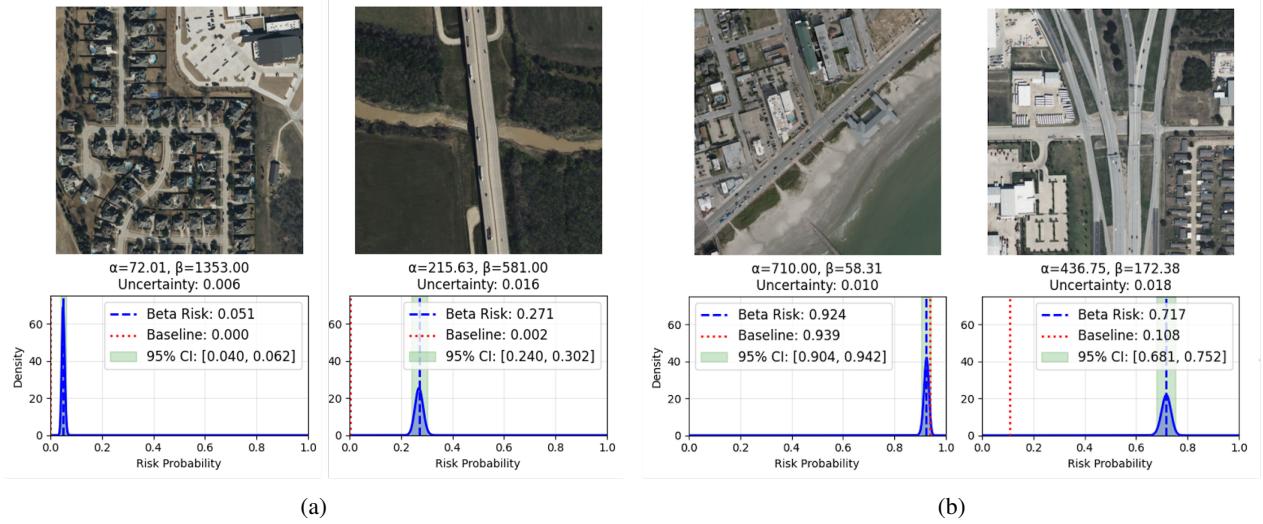


Figure 5: Qualitative Results for Unambiguous ("Easy") Cases (a: True Negatives, b: True Positives)

5.3 Visualizing Model Uncertainty

To be a trustworthy tool for risk assessment, a model must not only make accurate predictions but also provide a reliable measure of its own uncertainty. We visualize this uncertainty using a Beta distribution for each prediction. As shown in Figure 5, our model demonstrates well-calibrated confidence across a spectrum of cases, a crucial feature for real-world deployment.

For visually unambiguous locations, the model produces predictions with high confidence. For example, in a simple suburban neighborhood (Figure 5a, left), it predicts a low risk (0.051) with a correspondingly low uncertainty score (0.006), reflected in a sharp Beta distribution. Likewise, for a coastal road with high traffic density and high potential of distractions (Figure 5b, left), it correctly predicts a high risk (0.924) with high confidence (uncertainty of 0.010).

The model's utility is further demonstrated in more complex scenarios where it appropriately reduces its confidence. For a visually complex but safe highway overpass (Figure 5a, right), the model still correctly predicts low risk, but the wider Beta distribution indicates higher uncertainty. This nuanced confidence is also evident when assessing a complex highway interchange (Figure 5b, right); the model correctly predicts a high risk of 0.717, but acknowledges the significant uncertainty due to the challenging visual features.

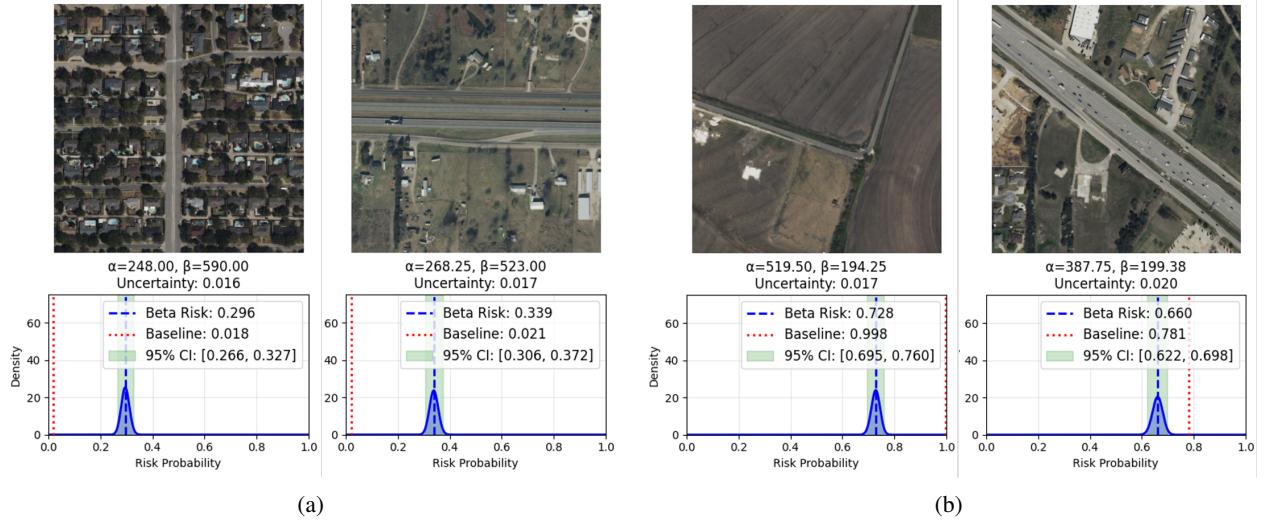


Figure 6: Interpreting Model Behavior on Ambiguous (“Hard”) Cases (a: False Negative, b: False Positive)

Crucially, the model’s rational expression of uncertainty extends to its failures, a characteristic vital for establishing trust. For false negatives (Figure 6a), where the model misses a crash, the low-risk predictions are consistently paired with wider, higher-uncertainty distributions. This correctly signals that the visual evidence was ambiguous, containing conflicting features (e.g., Figure 6a Left: an arterial road with many intersections within an otherwise low-risk residential area).

Similarly, for false positives (Figure 6b), the model flags locations as high-risk despite no recorded crashes, but again with reduced confidence. This behavior is highly interpretable, as the model correctly identifies latent risk factors, such as sharp (L-shape) turns or high-density highways. The prediction thus reflects a successful identification of hazardous features, while the increased uncertainty correctly marks them as borderline cases. This ability to temper certainty in response to visual complexity, especially when incorrect, distinguishes our model as a more reliable and interpretable system for risk assessment.

5.4 Case Study

To demonstrate the practical utility of our model, we conducted a case study of the San Antonio River Walk, a major tourist destination that presents a challenging environment with a complex mix of vehicular, pedestrian, and cyclist traffic. We generated risk predictions for over 140 locations in this area using Prob-MS and MSCM-MS.

The results, shown in Figure 7, highlight the superior performance of our approach. The baseline MSCM-MS model (middle panel) fails to identify close to half of the historical fatal crash locations (red diamonds), assigning them erroneously low risk scores. The baseline’s predictions also lack spatial coherence, exhibiting sharp, unrealistic gradients between adjacent points and producing polarized risk scores with few intermediate values.

In contrast, our Prob-MS model (right panel) correctly assigns elevated risk scores (yellow and orange) to a greater number of the known crash sites. This is exemplified at the intersection near Navarro St and Villita St, a known fatal crash location at the bottom-right in the map. While the baseline model misses this site, ours correctly assigns the area a high-risk score.

An analysis of the satellite and ground-level imagery reveals a confluence of latent risk factors not apparent from an overhead view alone. The location, a major entry point to the River Walk, is surrounded by numerous parking facilities. Ground-level images (Figure S1) show that these structures, combined with dense trees and building columns, create significant visual obstructions and blind spots for both drivers and pedestrians. This environment forces complex interactions: vehicles constantly enter and exit parking garages across wide pedestrian walkways as tourists navigate narrow sidewalks. Our model likely learned to associate this specific combination of visual clutter, unpredictable vehicle maneuvers, and high pedestrian-vehicle conflict with an elevated risk of a fatal crash.

Furthermore, our model generates a more nuanced and spatially coherent risk map where predictions transition smoothly across locations. This case study demonstrates that our model’s strong quantitative performance translates into more reliable, interpretable, and actionable safety assessments for complex urban environments.

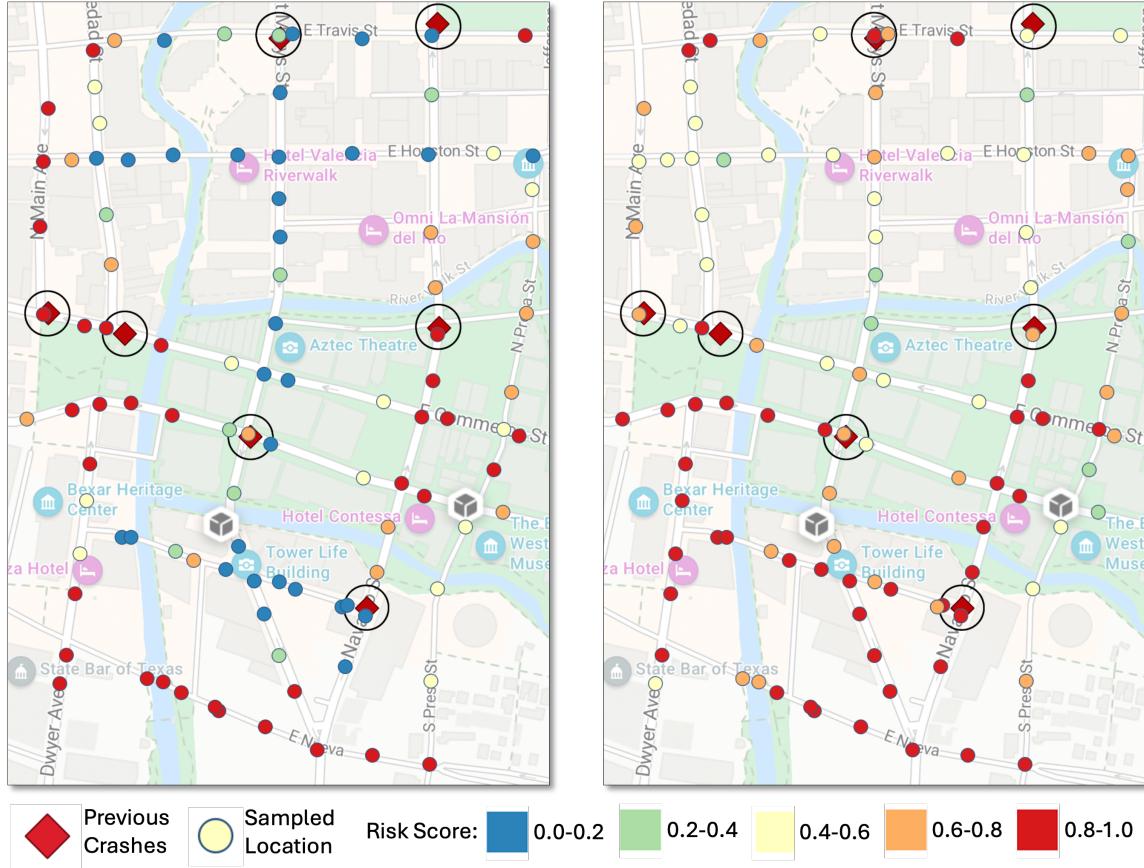


Figure 7: A case study of crash risk assessment for the San Antonio River Walk. Historical fatal crashes (red diamonds) serve as ground truth. (Left) The baseline MSCM-MS model exhibits low recall and spatially inconsistent predictions, with abrupt risk changes between adjacent points. (Right) Our Prob-MS model demonstrates superior recall by correctly identifying more crash sites and generates a more nuanced and spatially coherent risk field, providing a more realistic safety assessment. See the Result section for more.

6 Discussion and Conclusion

Our evaluation demonstrates that the proposed probabilistic framework yields a risk assessment model that is not only more effective but also more reliable and interpretable than deterministic baselines. By predicting a full Beta probability distribution instead of a single point-estimate, our model learns a more nuanced and less overconfident representation of risk. This trustworthiness is reinforced by its interpretable behavior; the model’s “mistakes” are often rational, such as flagging visually complex but historically safe highway interchanges as high-risk. This capacity to reason about visual factors and express nuanced confidence is highly valuable for practical applications, from enabling more sophisticated path planning for autonomous vehicles to allowing urban planners to confidently prioritize infrastructure improvements. Furthermore, by relying solely on publicly available satellite imagery, our method circumvents the significant privacy concerns associated with other data sources.

6.1 Ethical Considerations and Responsible Deployment

The ethical implications of deploying an AI tool for public safety are significant. As historical crash data may contain undiscovered biases, such as under-reporting in certain socioeconomic or geographic areas, a model used without critical oversight could perpetuate inequities. We therefore emphasize that this model is designed as a decision-support tool to augment, not replace, human expertise.

A key feature for responsible, human-in-the-loop deployment is the model’s ability to signal its own uncertainty, which can serve as a bias and fairness mitigation tool. High uncertainty in any prediction (whether high-risk or low-risk) can flag regions with potential data disparities or under-reporting. For instance, a visually complex area with high

uncertainty and a low-risk prediction may indicate a dangerous false negative due to a lack of historical crash data. These uncertain predictions should act as a flag for human experts to conduct a more detailed investigation, enabling a more equitable allocation of safety resources.

6.2 Limitations

This study has several limitations that open avenues for future research. Our model estimates static risk and does not account for dynamic variables like real-time traffic or weather; future work should focus on integrating these data streams. Our study is also geographically constrained to Texas, and validation on diverse international datasets is a critical next step to ensure generalizability. Furthermore, this work can be extended by exploring a learned weighting mechanism for the centrality and size components of our procedural labeling scheme. Finally, while our model identifies strong correlations, future work could explore methods for moving toward causal inference.

In conclusion, this work demonstrates that moving from deterministic point-estimates to a full probabilistic framework is a crucial step toward creating more reliable and trustworthy AI for public safety. By learning to predict a Beta probability distribution from satellite imagery, our model not only outperforms existing baselines in identifying high-risk locations but also provides the well-calibrated uncertainty estimates that are vital for interpretable, human-in-the-loop decision-making in applications from urban planning to autonomous navigation.

7 Conclusion

This work presents a deep learning framework for reliable roadway risk assessment that quantifies uncertainty. Instead of a single risk score, our model predicts a full Beta probability distribution to provide a more comprehensive hazard assessment. This is achieved using a procedural labeling technique with data augmentation to supervise uncertainty, and a compound loss function that jointly optimizes for classification accuracy and probabilistic calibration.

Our model significantly outperforms existing baselines, with a 17-23% relative improvement in recall on high-risk locations, up to 17% in ECE on calibration, and about 11% more stable. More importantly, it yields interpretable predictions, reliably signaling its uncertainty in ambiguous cases. By delivering a more robust and trustworthy assessment of roadway risk, our work represents a crucial step toward the responsible deployment of AI in high-stakes applications, such as public safety, urban planning, and autonomous navigation.

Acknowledgement

This material is partially based upon work supported by the National Science Foundation under 2401860 and 2526487. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and the funders have no role in the study design, data collection, analysis, or preparation of this article.

Portions of this research were conducted with the advanced computing resources provided by the High Performance Computing Research Center at Texas A&M University-San Antonio.

References

- [1] WHO, “Road traffic injuries,” *World Health Organization*, 2023, accessed: 2025-05-22. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/road-traffic-injuries>
- [2] ———, “Global status report on road safety,” *World Health Organization*, 2018, accessed: 2025-05-22. [Online]. Available: <https://www.who.int/publications/item/9789241565684>
- [3] CISA, “Critical infrastructure sectors,” *Cybersecurity and Infrastructure Security Agency*, 2024, accessed: 2024-03-22. [Online]. Available: <https://www.cisa.gov/topics/critical-infrastructure-security-and-resilience/critical-infrastructure-sectors>
- [4] C. Caliendo *et al.*, *Accident Analysis & Prevention*, vol. 39, no. 4, pp. 657–670, 2007.
- [5] J. Tamerius, X. Zhou, R. Mantilla, and T. Greenfield-Huitt, “Precipitation effects on motor vehicle crashes vary by space, time, and environmental conditions,” *Weather, Climate, and Society*, vol. 8, no. 4, pp. 399–407, 2016.
- [6] C. Zhu, B. Dadashova, C. Lee, X. Ye, and C. T. Brown, “Equity in non-motorist safety: Exploring two pathways in houston,” *Transportation research part D: transport and environment*, vol. 132, p. 104239, 2024.
- [7] B. G. Simons-Morton, F. Guo, S. G. Klauer, J. P. Ehsani, and A. K. Pradhan, “Keep your eyes on the road: Young driver crash risk increases according to duration of distraction,” *Journal of Adolescent Health*, vol. 54, no. 5, pp. S61–S67, 2014.

- [8] A. Pembuain *et al.*, “The effect of road infrastructure on traffic accidents,” in *11th Asia Pacific Transportation and the Environment Conference (APTE 2018)*. Atlantis Press, 2019, pp. 176–182.
- [9] T. Huang *et al.*, “Highway crash detection and risk estimation using deep learning,” *Accident Analysis & Prevention*, vol. 135, p. 105392, 2020.
- [10] D. Jaroszwecki and T. McNamara, “The influence of rainfall on road accidents in urban areas: A weather radar approach,” *Travel behaviour and society*, vol. 1, no. 1, pp. 15–21, 2014.
- [11] C. Gu, J. Xu, C. Gao, M. Mu, G. E, and Y. Ma, “Multivariate analysis of roadway multi-fatalities crashes using association rules mining and rules graph structures: A case study in china,” *Plos one*, vol. 17, no. 10, p. e0276817, 2022.
- [12] C. Carrodano, “Data-driven risk analysis of nonlinear factor interactions in road safety using bayesian networks,” *Scientific Reports*, vol. 14, no. 1, p. 18948, 2024.
- [13] I. Ahmed, “Road infrastructure and road safety,” *Transport and Communications Bulletin for Asia and the Pacific*, vol. 83, pp. 19–25, 2013.
- [14] W. Song, S. Workman, A. Hadzic, X. Zhang, E. Green, M. Chen, R. Souleyrette, and N. Jacobs, “Farsa: Fully automated roadway safety assessment,” in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2018, pp. 521–529.
- [15] G. Cheng, R. Cheng, S. Zhang, and X. Sun, “Risk evaluation method for highway roadside accidents,” *Advances in Mechanical Engineering*, vol. 11, no. 1, p. 1687814018821743, 2019.
- [16] Q. Ma, H. Yang, Z. Wang, K. Xie, and D. Yang, “Modeling crash risk of horizontal curves using large-scale auto-extracted roadway geometry data,” *Accident Analysis & Prevention*, vol. 144, p. 105669, 2020.
- [17] Y.-J. Joo *et al.*, “A generalized driving risk assessment on high-speed highways using field theory,” *Analytic Methods in Accident Research*, vol. 40, p. 100303, 2023.
- [18] V. de Almeida Guimarães *et al.*, “Evaluating the sustainability of urban passenger transportation by monte carlo simulation,” *Renewable and Sustainable Energy Reviews*, vol. 93, pp. 732–752, 2018.
- [19] L. Al-Sharif *et al.*, “The use of monte carlo simulation in evaluating the elevator round trip time under up-peak traffic conditions and conventional group control,” *Building Services Engineering Research and Technology*, vol. 33, no. 3, pp. 319–338, 2012.
- [20] S. Jeon and B. Hong, “Monte carlo simulation-based traffic speed forecasting using historical big data,” *Future generation computer systems*, vol. 65, pp. 182–195, 2016.
- [21] S. Moosavi, M. H. Samavatian, S. Parthasarathy, R. Teodorescu, and R. Ramnath, “Accident risk prediction based on heterogeneous sparse data: New dataset and insights,” in *Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2019, pp. 33–42.
- [22] S. He, M. A. Sadeghi, S. Chawla, M. Alizadeh, H. Balakrishnan, and S. Madden, “Inferring high-resolution traffic accident risk maps based on satellite imagery and gps trajectories,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 11977–11985.
- [23] A. Najjar *et al.*, “Combining satellite imagery and open data to map road safety,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.
- [24] Y. Li *et al.*, “Label bias: A pervasive and invisibilized problem,” *Notices of the American Mathematical Society*, vol. 71, no. 8, pp. 1069–1077, 2024.
- [25] C. Chen and S. S. Sundar, “Is this ai trained on credible data? the effects of labeling quality and performance bias on user trust,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 2023, pp. 1–11.
- [26] X. Xing, G. Liang, C. Wang, N. Jacobs, and A.-L. Lin, “Self-supervised learning application on covid-19 chest x-ray image classification using masked autoencoder,” *Bioengineering*, vol. 10, no. 8, p. 901, 2023.
- [27] L. Liu, Y. Wang, J. Chang, P. Zhang, G. Liang, and H. Zhang, “Llrhnet: multiple lesions segmentation using local-long range features,” *Frontiers in Neuroinformatics*, vol. 16, p. 859973, 2022.
- [28] J. Zulu, B. Han, I. Alsmadi, and G. Liang, “Enhancing machine learning based sql injection detection using contextualized word embedding,” in *Proceedings of the 2024 ACM Southeast Conference*, 2024, pp. 211–216.
- [29] R. Jonnala, G. Liang, J. Yang, and I. Alsmadi, “Exploring the potential of large language models in public transportation: San antonio case study,” 2025.

- [30] S.-C. Lin, Y. Su, G. Liang, Y. Zhang, N. Jacobs, and Y. Zhang, “Estimating cluster masses from sdss multiband images with transfer learning,” *Monthly Notices of the Royal Astronomical Society*, vol. 512, no. 3, pp. 3885–3894, 2022.
- [31] G. Pereyra, G. Tucker, J. Chorowski, Ł. Kaiser, and G. Hinton, “Regularizing neural networks by penalizing confident output distributions,” in *International Conference on Learning Representations*, 2017.
- [32] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, “On calibration of modern neural networks,” in *International conference on machine learning*. PMLR, 2017, pp. 1321–1330.
- [33] G. E. Hinton *et al.*, “Distilling the knowledge in a neural network,” *ArXiv*, vol. abs/1503.02531, 2015.
- [34] A. Kumar *et al.*, “Trainable calibration measures for neural networks from kernel mean embeddings,” in *International Conference on Machine Learning*, 2018, pp. 2805–2814.
- [35] G. Liang, Y. Zhang, X. Wang, and N. Jacobs, “Improved trainable calibration method for neural networks on medical imaging classification,” in *British Machine Vision Conference (BMVC)*, 2020.
- [36] M. Chidambaram and R. Ge, “On the limitations of temperature scaling for distributions with overlaps,” in *International Conference on Learning Representations*, 2023.
- [37] L. N. Vaserstein, “Markov processes over denumerable products of spaces, describing large systems of automata,” *Problemy Peredachi Informatsii*, vol. 5, no. 3, pp. 64–72, 1969.
- [38] I. Csiszár, “I-divergence geometry of probability distributions and minimization problems,” *The annals of probability*, pp. 146–158, 1975.
- [39] H. Cramér, “On the composition of elementary errors,” *Scandinavian Actuarial Journal*, vol. 1928, no. 1, pp. 13–74, 1928.
- [40] G. Liang, J. Zulu, X. Xing, and N. Jacobs, “Unveiling roadway hazards: Enhancing fatal crash risk estimation through multiscale satellite imagery and self-supervised cross-matching,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 17, pp. 535–546, 2024.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [42] A. Krizhevsky *et al.*, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, 2012.
- [43] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [44] ——, “Sgdr: Stochastic gradient descent with warm restarts,” *arXiv preprint arXiv:1608.03983*, 2016.

Supplementary Materials

San Antonio River Walk Ground-Level Imagery

Figure S1 shows the ground-level image at 146-Navarro-St, San Antonio, TX, USA, one main entry point to the San Antonio River Walk area. A previous fatal crash also occurred at this location (i.e., the bottom-right one in Figure 7).

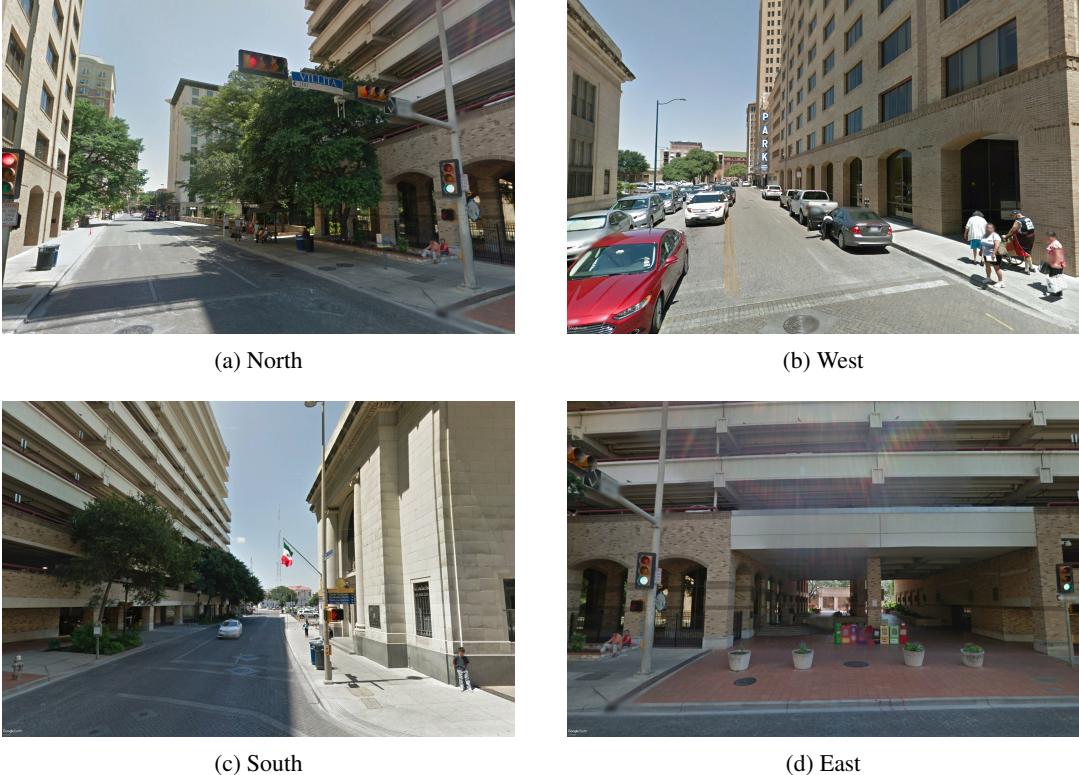


Figure S1: Ground-level images for the four directions at 146-Navarro-St, San Antonio, TX, USA, one main entry point to the San Antonio River Walk area.

Dataset

This study utilizes the comprehensive, multi-scale satellite imagery dataset provided by MSCM [40]. The dataset covers diverse regions in Texas, USA, including the Gulf Coast, Hill Country, and Prairies and Lakes regions.

The dataset contains a total of 240,828 satellite images. The images for each location are provided at three distinct levels of detail, each with a resolution of 768×768 pixels. The images for each location are provided at three distinct levels of detail: 1.1943 m/pixel, 0.5972 m/pixel, and 0.2986 m/pixel. Examples of this multi-scale imagery are shown in Figure S2.

The data is sampled from 80,276 distinct locations, which are categorized into positive and negative classes. The *positive* class consists of 16,451 locations where at least one fatal crash occurred between 2010 and 2020; of these, 1,185 locations experienced multiple fatal crashes within a 50-meter radius. The remaining locations serve as the *negative* class, having no recorded fatal crashes through the end of 2020. These negative samples were selected using specific criteria, ensuring they were within 1250 meters of a fatal crash location but at least 250 meters away from any such site. To create a challenging learning environment, approximately 70% of the negative samples were designated as hard negatives by sampling them along primary and secondary roads. The other 30% were sampled randomly to represent a broader variety of environments, including open spaces.



Figure S2: Multi-Scale Satellite Imagery Inputs with Three Detail Levels. From left to right: 1.1943 m/pixel, 0.5972 m/pixel, and 0.2986 m/pixel.

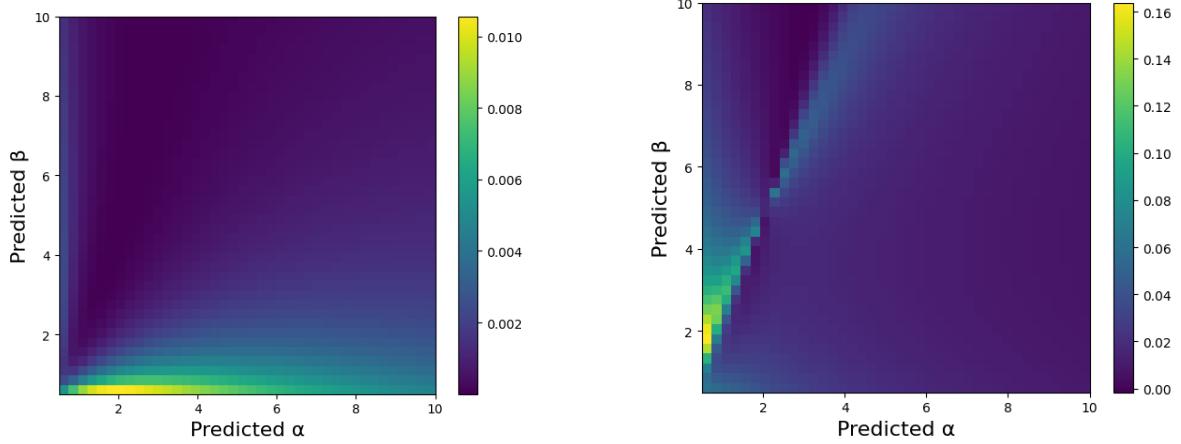


Figure S3: Comparison between true and surrogate Wasserstein-2 distance for Beta distributions. **Left:** Absolute Error. **Right:** Relative Error.

Wasserstein-2 Surrogate Analysis

To evaluate the accuracy of our surrogate Wasserstein-2 loss, we computed the true squared Wasserstein-2 distance between a fixed target distribution $\text{Beta}(2, 5)$ and a range of predicted Beta distributions with $\alpha, \beta \in [0.5, 10]$. For each pair of predicted parameters, the true distance was estimated via numerical integration of the quantile functions, while the surrogate distance was computed using the closed-form mean–variance expression defined in our loss. The resulting absolute and relative differences are visualized in Figure S3. Both plots demonstrate that the surrogate closely matches the true distance, with errors typically on the order of 10^{-3} to 10^{-2} and only slightly increasing in extreme parameter regions.

Implementation Details

All models are built upon a ResNet-50 backbone and trained for 75 epochs using the AdamW [43] optimizer with a CosineAnnealingWarmRestarts [44] learning rate schedule. The initial learning rate for the backbone and classifier was set to $1e^{-4}$, while the distribution learning head used a rate of 0.02. The batch size was 128 for single-scale models and 48 for multi-scale models. Based on our hyperparameter analysis, the final weights for our compound loss function were set to $\lambda_1 = 5$ and $\lambda_2 = 1$ to prioritize recall for this safety-critical task. The same data augmentation pipeline are used for all models. The random seeds are set to 0. All experiments were conducted on two NVIDIA A100 GPUs. The rest of this section, provided the complete list of hyperparameters and data augmentation settings.

Model Architecture

- **Backbone:** All models use a ResNet-50 architecture with weights pre-trained on ImageNet.
- **Modification:** A 1x1 convolutional layer was inserted before the final Global Average Pooling .

MSCM Pre-training and Pre-trained Checkpoint Selection

The MSCM weights, used to initialize our proposed model and the corresponding baselines, were generated by following the pre-training procedure described in the original work.

- **Setup:** The pre-training used a contrastive learning approach with an InfoNCE and classification loss. It was run for 25 epochs with a batch size of 64, using the AdamW optimizer with a learning rate of 10^{-3} and a CosineAnnealingWarmRestarts schedule ($T_0=10$, $T_mult=2$).
- **Checkpoint Selection:** A model checkpoint was saved after each pre-training epoch. To select the optimal checkpoint, each of the 25 checkpoints was used to initialize a single-scale classification model, which was then fine-tuned for 5 epochs on the downstream task (batch size 128, AdamW, learning rate 10^{-4} , CosineAnnealingWarmRestarts schedule). The checkpoint that yielded the best performance after this short fine-tuning process was selected for all main experiments.

Main Training for Proposed Probabilistic Models

- **Optimizer:** AdamW.
- **Learning Rates (LR):** We used different learning rates for distinct parts of the model:
 - Feature Extractor Backbone: 10^{-4}
 - Beta Distribution Learning Head: 0.02
 - Auxiliary Classification Head: 10^{-4}
- **LR Schedule:** CosineAnnealingWarmRestarts with scheduler parameters $T_0=10$ and $T_mult=2$.
- **Epochs:** All models were trained for 75 epochs. The epoch with the highest accuracy on the test set is chosen as the best model.
- **Batch Size:** 128 for single-scale models; 48 for multi-scale models.

Loss Function and Hyperparameters

- **Compound Loss Weights:** The reported results use $\lambda_1 = 5$ (for BCE loss) and $\lambda_2 = 1$ (for Wasserstein loss).
- **BCE Loss:** Class imbalance was handled using inverse frequency weights applied to the BCE loss: [1.25948, 4.85382].
- **Procedural Beta Distribution Parameters:**
 - base_K: 22.0
 - ϵ : 0.08
 - min_positive_risk_mean: 0.18
 - min_concentration_positives: 18.0
 - Influence Score Weights: weight_distance=0.7, weight_crop_size=0.3.

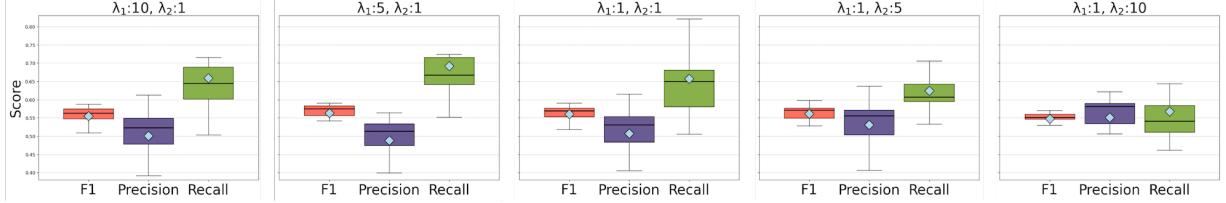


Figure S4: Hyperparameter Analysis: Precision-Recall Trade-off

Table S1: Ablation Study on Loss Weights

λ_1	λ_2	F1 Score (\uparrow)	Precision (\uparrow)	Recall (\uparrow)
10	1	0.5880	0.5358	0.6514
5	1	0.5981	0.5607	0.6409
1	1	0.5905	0.5406	0.6505
1	5	0.5979	0.5932	0.6026
1	10	0.5855	0.5577	0.6163

Data Augmentation

The augmentation pipelines were used for our proposed models versus the baselines.

- **For Proposed Probabilistic Models:**
 - Random Crop Ratio: (0.5, 1.0)
 - Random horizontal flip ($p=0.5$)
 - Random vertical flip ($p=0.5$)
 - Random rotations (from -90 to 90 degrees)
 - ColorJitter (brightness/contrast/saturation: [0.6, 1.4], hue: [0.0, 0.1])
- **For Baseline Models:**
 - Random Crop Ratio: (0.3, 1.0)
 - Random horizontal flip ($p=0.5$)
 - Random vertical flip ($p=0.5$)
 - Random rotations (from -90 to 90 degrees)
 - ColorJitter (brightness/contrast/saturation: [0.6, 1.4], hue: [0.0, 0.1])

Hardware

- All experiments were conducted on two NVIDIA A100 GPUs, each with 40GB of memory.

Hyperparameter Analysis: Effect of Loss Weights

To understand how the components of our compound loss function influence model behavior, we conducted a hyperparameter analysis on the loss weights, λ_1 (for the classification loss, \mathcal{L}_{BCE}) and λ_2 (for the Beta distribution loss, $\mathcal{L}_{W_2^2}$). Table S1 demonstrates that these weights serve as a practical lever to tune the model’s predictive trade-offs for different application needs. See the supplementary materials for a detailed analysis.

Our analysis confirms two key trends. First, increasing the weight of the classification loss (λ_1) makes the model prioritize **Recall**. As shown in Table S1, increasing λ_1 to 10 yields the highest recall score of 0.6514. This indicates that a stronger emphasis on the classification task pushes the model to more aggressively identify all potential high-risk locations, which is critical for safety applications.

Conversely, increasing the weight of the Beta distribution loss (λ_2) encourages a more **balanced and precise** model. The box plots in Figure S4, which show the performance distribution over 25 epochs, provide additional insight. They visually confirm that as λ_2 increases, the median precision rises while recall moderately decreases, bringing the two metrics into closer alignment. This is exemplified by the $\lambda_1 = 1, \lambda_2 = 5$ configuration, which achieves the highest precision of all tested settings (0.5932) while maintaining a strong recall (0.6026), as detailed in Table S1. This

demonstrates that a stronger emphasis on the distribution-matching loss component encourages a more conservative model that makes fewer, but more accurate, high-risk predictions.

This analysis provides clear guidance for hyperparameter selection based on the desired outcome. For a safety-critical system where failing to identify a hazard is the worst-case scenario, a higher λ_1 is optimal. For applications requiring high confidence in positive predictions to efficiently allocate resources, a higher λ_2 would be chosen. For the main results reported in this paper, we selected the $\lambda_1 = 5, \lambda_2 = 1$ configuration, as it achieved the highest F1-score and maintained a strong recall, offering an excellent balance for our primary task.