

Benchmarking YOLOv8-YOLOv12 Models for Image-Space Near-Miss Detection in Construction Safety

1 **Shahid Ali¹, Ammar Alzarrad², Sudipta Chowdhury³, Husnu S. Narman⁴**

2 ¹Department of Computer Sciences and Electrical Engineering, Marshall University, One John
3 Marshall Drive, Huntington, WV 25755; e-mail: shahidali@marshall.edu

4 ²Department of Civil Engineering, Marshall University, One John Marshall Drive, WV 25755; E-
5 mail: alzarrad@marshall.edu

6 ³Department of Mechanical and Industrial Engineering, Marshall University, One John Marshall
7 Drive, Huntington, WV 25755; e-mail: chowdhurys@marshall.edu

8 ⁴Department of Computer Sciences and Electrical Engineering, Marshall University, One John
9 Marshall Drive, Huntington, WV 25755; e-mail: narman@marshall.edu

10 *** Correspondence:**

11 Shahid Ali
12 shahidali@marshall.edu

13 **Keywords: Near-miss detection, construction safety monitoring, YOLO object detection,**
14 **worker-equipment interaction, computer vision in construction.**

15 **Abstract**

16 Automated near-miss detection can support proactive construction safety monitoring by identifying
17 hazardous worker-equipment interactions before injuries occur. This study presents a systematic
18 evaluation framework for comparing 24 publicly available pretrained YOLO object detection models,
19 spanning YOLOv8 through YOLOv12, for image-based near-miss detection in construction
20 environments. The evaluation used 38 publicly available construction-related videos, standardized to
21 1920×1080 resolution and sampled at 1 frame per second, resulting in approximately 903 evaluation
22 frames across varied work-zone conditions, camera perspectives, object scales, clutter, motion blur,
23 and partial occlusion. A rule-based pipeline identified person-equipment interactions using centroid-
24 based image-space proximity under a primary 250-pixel threshold. Because complete camera
25 calibration metadata were unavailable, this study treats this threshold as an operational image-space
26 criterion rather than a universally calibrated physical distance. This study evaluates object-detection
27 performance using precision, recall, and F1 score against a human-corrected consensus ground truth
28 derived from YOLOv8x candidate annotations and independently reviewed by three domain experts.
29 The analysis assesses near-miss performance using frame-level recall, spatial overlap, fuzzy centroid
30 matching, and threshold-sensitivity analysis. Results showed that YOLOv8x, YOLOv8l, and
31 YOLOv9c achieved the strongest near-miss performance. The findings highlight the importance of
32 stable frame-level detection and provide an exploratory, reproducible benchmark of comparative
33 YOLO model behavior under selected image-space construction-video conditions.

34 **1 Introduction**

35 Construction projects involve constantly changing work conditions, moving equipment, shifting site
36 layouts, and frequent interactions between workers and machinery. These conditions make
37 construction sites particularly vulnerable to safety risks, including near-miss incidents where workers
38 narrowly avoid injury. According to the U.S. Bureau of Labor Statistics (BLS, 2026), there were 5,070
39 fatal work injuries in the United States in 2024, down 4.0% from 5,283 in 2023. The fatal work injury
40 rate was 3.3 fatalities per 100,000 full-time equivalent workers in 2024, compared with 3.5 in 2023.
41 These national figures highlight the continued need for proactive safety monitoring in high-risk
42 industries such as construction, where proximity hazards and improper use of personal protective
43 equipment remain persistent concerns (Olimat et al., 2025). Traditional safety monitoring, which often
44 relies on manual inspections, is labor-intensive, subjective, and limited in its ability to detect hazards
45 in real time (López et al., 2025; Lim et al., 2025). These limitations become more pronounced as project
46 complexity and pace increase, particularly when brief worker-equipment interactions must be
47 identified quickly (López et al., 2025).

48 Advances in computer vision and deep learning offer new opportunities for proactive safety monitoring
49 and automated assessment of civil infrastructure. Recent studies have demonstrated the use of drones,
50 depth sensors, and machine-learning models for roof-condition inspection, solar-panel site assessment,
51 earthquake-damage assessment, railroad gauge measurement, railway bolt detection, and construction-
52 site firearm/tool/person detection (Alzarrad et al., 2021; Alzarrad et al., 2022; Kizilay et al., 2024; Le
53 et al., 2025; Alzarrad et al., 2025; Damai et al., 2025). The YOLO (You Only Look Once) family of
54 object detectors is widely used in real-time applications because it balances detection speed and
55 accuracy (Lim et al., 2025; Alotaibi and Ma, 2025). Prior studies have applied YOLO models to PPE
56 compliance (López et al., 2025), behavior recognition (Alotaibi and Ma, 2025), drowsiness detection
57 (Onososen et al., 2025), and risk-zone alerts (Dzeng et al., 2025). However, most existing studies focus
58 on a single YOLO variant or a specific detection task, leaving limited evidence on how newer YOLO
59 versions and model sizes perform in near-miss scenarios where spatial relationships between workers
60 and equipment are central.

61 This study introduces a comparative framework to evaluate YOLOv8 through YOLOv12 for near-miss
62 detection in construction environments. Using centroid-based proximity analysis, the proposed
63 pipeline processes frame-level video data to identify worker-equipment interactions. Outputs include
64 annotated images, event logs, and detection metrics such as F1 score, recall, and spatial overlap,
65 allowing cross-version comparison of detection behavior and risk sensitivity. This study seeks to
66 answer the following questions: (1) How do different YOLO versions and sizes compare in detecting
67 near-misses? (2) How do YOLO model generation and model size relate to detection accuracy in spatial
68 risk contexts? The study provides a reproducible evaluation framework for construction safety research
69 by standardizing inputs and evaluation criteria. The remainder of this paper is organized as follows.
70 The Literature Review summarizes prior YOLO-based safety-monitoring studies and identifies gaps
71 in interaction-level near-miss evaluation. The Methodology section describes the video dataset, the
72 model-selection process, the annotation protocol, the proximity-based near-miss logic, and the
73 evaluation metrics. The Results section presents object-detection and near-miss performance across all
74 24 YOLO variants, including threshold-sensitivity analysis. The Discussion interprets the findings in
75 relation to model size, detection consistency, deployment constraints, and false-positive/false-negative
76 tradeoffs. The Conclusion summarizes the study's contributions, practical implications, and directions
77 for future research.

78 **2 Literature Review**

79 This section reviews recent YOLO-based and computer vision approaches to construction safety
80 monitoring, with emphasis on studies published within the last five years. The review focuses on three
81 areas relevant to this study: object-level safety detection, interaction-level worker-equipment risk
82 assessment, and methodological issues related to ground-truth validation and proximity-based near-
83 miss definitions. Near-miss detection differs from conventional object detection because it requires
84 evaluating spatial relationships between workers and equipment rather than identifying isolated objects
85 alone. YOLO-based models play a central role in safety monitoring tasks such as PPE compliance,
86 drowsiness detection, and collision alerts, with many studies reporting strong detection performance
87 (Onososen et al., 2025; Mohy et al., 2025; Dzeng et al., 2025). However, these studies primarily address
88 static PPE compliance, isolated unsafe behaviors, object-level detection, or single-model
89 configurations. They rarely evaluate interaction-level risk, where the spatial relationship between
90 workers and equipment determines whether a detection becomes safety-critical. Prior work highlights
91 the impact of occlusion, background clutter, lighting variation, small-object detection, and edge-
92 computing constraints on real-world deployment; however, these factors remain seldom examined
93 within a standardized multi-model near-miss evaluation framework. Although prior research proposes
94 proximity-based frameworks for near-miss detection (Lim et al., 2025), existing studies remain limited
95 to specific domains or individual model configurations, and no study provides a systematic evaluation
96 across multiple YOLO generations using real-world construction-related footage. Existing YOLO-
97 based methods, therefore, still lack a standardized evaluation framework for interaction-level risk
98 assessment in near-miss detection scenarios.

99 Over the past five years, computer-vision research has expanded the role of YOLO-based models in
100 safety monitoring beyond static object recognition. Earlier video-surveillance studies demonstrated
101 that YOLO detections can be integrated with tracking and data-association methods to improve
102 robustness under scale variation, occlusion, and motion uncertainty (Ait Abdelali et al., 2021).
103 Subsequent work has explored hybrid CNN-Transformer detectors, real-scene safety datasets, and
104 data-augmentation strategies to address small-object detection, class imbalance, and environmental
105 variability in complex scenes (Guo et al., 2022; Yongxiang et al., 2022; Jiang et al., 2024). In
106 construction-specific applications, recent studies have applied YOLO-based models to PPE detection,
107 unsafe-behavior recognition, site-risk assessment, and UAV-based monitoring, showing strong
108 promise while also reporting persistent challenges related to lighting variation, occlusion, cross-site
109 generalization, and edge-device deployment (Feng et al., 2024; Zou and Hu, 2024; Lyu et al., 2025;
110 Zhou et al., 2025). Together, these studies show that the field has progressed from basic object
111 detection toward more context-aware, deployment-oriented monitoring; however, most prior work still
112 evaluates detection accuracy at the object or task level rather than assessing how model outputs
113 translate into near-miss identification at the interaction level.

114 Beyond construction PPE and near-miss applications, related civil infrastructure studies further
115 demonstrate the value of AI-enabled visual sensing for safety-critical assessment. Drone-based and
116 neural network approaches support the evaluation of roof suitability for solar panel installation and the
117 assessment of roof conditions, reducing reliance on manual inspection in elevated or difficult-to-access
118 environments (Alzarrad et al., 2021; Alzarrad et al., 2022). Similarly, drone imagery and fine-tuned
119 deep-learning models have been applied to earthquake-damage assessment, while depth-sensor and
120 machine-learning approaches have supported railroad gauge measurement and missing-bolt detection
121 for railway safety applications (Kizilay et al., 2024; Le et al., 2025; Damai et al., 2025). These studies
122 reinforce the broader relevance of computer vision for automating visual inspection and safety
123 monitoring. However, they do not directly address near-miss detection at the interaction level for
124 worker-equipment interactions in dynamic construction scenes.

125 Concerns regarding ground-truth validity and evaluation bias further highlight limitations in existing
126 studies. While advanced detection frameworks achieve strong performance in safety monitoring tasks
127 (Lee et al., 2023; Alotaibi and Ma, 2025), most studies fail to report inter-annotator agreement, expert
128 labeling qualifications, and mechanisms to mitigate circular validation in model-assisted annotations.
129 A broader review of construction safety research indicates that most AI-driven monitoring studies omit
130 transparent human validation protocols (Olimat et al., 2025). Moreover, prior reviews of occlusion-
131 handling approaches show that object detectors often degrade under partial occlusion, cluttered
132 backgrounds, and complex surveillance conditions, and these limitations are not fully captured by IoU-
133 based metrics alone (Ouardirhi et al., 2024). These findings underscore the necessity of independently
134 human-validated datasets and documented annotation reliability when evaluating near-miss detection
135 in safety-critical environments.

136 Methodological inconsistency is also evident in the calibration of proximity thresholds used to define
137 near-miss events. Prior studies acknowledge that pixel-based distance thresholds are highly sensitive
138 to camera geometry, field of view, lens distortion, and mounting configuration, yet such parameters
139 are rarely reported in sufficient detail (Lim et al., 2025; Maulana and Hardiansyah, 2025). The literature
140 suggests that robust near-miss evaluation should combine sensitivity analysis across multiple
141 thresholds with camera-aware calibration to establish meaningful correspondence between pixel
142 distance and physical separation. Treating proximity thresholds as tunable hyperparameters aligns with
143 best practices in spatial risk assessment, though explicit calibration documentation remains essential
144 for reproducibility and cross-site generalization.

145 Addressing these gaps, the present study provides a comparative evaluation of 24 YOLO variants
146 spanning YOLOv8-YOLOv12 for near-miss detection using publicly available construction-related
147 video footage. By shifting the emphasis from isolated object-detection metrics to frame-level recall
148 and interaction-level risk identification, this work establishes a reproducible benchmarking framework
149 tailored to safety-critical applications. The findings demonstrate that consistent temporal detection and
150 contextual spatial reasoning are more influential for near-miss identification than absolute proximity
151 distance alone, offering practical insights beyond conventional object-detection evaluation paradigms.
152 Future research should extend this framework through expanded human-validated datasets, camera-
153 geometry-based calibration, temporal modeling, edge-deployment evaluation, and multimodal sensing
154 to further enhance robustness under occlusion and dynamic site conditions. Table 1 summarizes
155 representative studies from 2021 to 2025 across video surveillance, construction safety, UAV
156 monitoring, PPE detection, occlusion handling, and proximity-based near-miss detection.

157 **Table 1. Comparative Analysis of YOLO-Based Approaches for Interaction-Level Risk**
158 **Assessment**

Study	Domain	Method/Focus	Performance	Key Limitation
Ait Abdelali et al. (2021)	Video surveillance / traffic	YOLO + Kalman filtering + Hungarian tracking	94.10% detection accuracy; 92.50% tracking accuracy	Traffic domain; not construction near-miss detection

Guo et al. (2022)	Industrial object detection	MSFT-YOLO / hybrid YOLO-Transformer detection	Improved defect-detection performance	Industrial defect domain; transfer to construction safety untested
Yongxiang et al. (2022)	Real-scene object detection	Real-scene detection dataset for object detection	Dataset contribution	Limited safety semantics and no near-miss interaction labels
Lee et al. (2023)	Construction	YOLACT + DeepSORT + PPE	91.3% accuracy	No inter-annotator agreement reporting
Feng et al. (2024)	Construction safety	YOLO-based PPE / unsafe-behavior or site-risk detection	Reported improved construction-site safety detection	Limited interaction-level near-miss benchmarking
Jiang et al. (2024)	Construction / PPE detection	Data augmentation for robust PPE recognition	Improved robustness under limited data	Rare near-miss and worker–equipment interaction events underrepresented
Zou and Hu (2024)	UAV / small-object hazard detection	YOLO-based aerial small-object detection	Improved small-object detection	Precision–recall tradeoff under tiny targets and changing viewpoints
Onososen et al. (2025)	Construction	YOLOv8 + drowsiness detection	92% mAP	Single static behavior, not interactions
Alzarrad et al. (2025)	Construction safety	AI-driven detection of people, tools, and firearms on construction sites	Reported automated multi-class detection performance	Focused on object-level threat detection; not worker-equipment near-miss interaction benchmarking
Mohy et al. (2025)	Construction	Computer vision model for automated safety compliance	Reported automated safety-compliance detection performance	Focused on safety compliance; not interaction-level near-miss benchmarking

Dzeng et al. (2025)	Construction	YOLOv7 + dynamic collision alert	0.998 mAP	Simulated virtual environment, not real footage
Lim et al. (2025)	Traffic / proximity	YOLOv7CNeB + distance indicators	0.997 mAP	Single model, traffic domain only
Alotaibi and Ma (2025)	Construction	Vision Transformer + unsafe-worker behavior detection	93.2% precision	No ground-truth validation methodology
Ouardirhi et al. (2024)	Smart video surveillance / object detection	Survey of occlusion-handling approaches	Review study	Addresses occlusion in surveillance but not construction-specific near-miss spatial context
Lyu et al. (2025)	UAV / edge detection	Lightweight YOLO for aerial edge deployment	Reported improved edge-oriented detection	Cross-site domain shift and edge-performance tradeoffs remain
Zhou et al. (2025)	UAV / small-object detection	SMA-YOLO with multi-scale attention and feature fusion	Improved small-object detection performance	Small-object recall, occlusion, and deployment tradeoffs remain
Saeedizadeh et al. (2025)	Autonomous driving / image-based detection	Survey of deep-learning object-detection methods	Review of 90+ models	Deployment robustness under occlusion, low light, and complex scenes remains challenging
Olimat et al. (2025)	Construction safety / bibliometric review	Data-driven review of construction safety research	14,174 publications reviewed; 61.8% recent studies from 2016–2025	Most studies lack transparent validation

160 This section describes the dataset, model selection process, near-miss detection logic, ground truth
161 annotation protocol, and the metrics used to evaluate both object detection and interaction-level
162 performance across 24 YOLO variants. This study introduces a systematic framework for evaluating
163 the performance of state-of-the-art object detection models, YOLOv8 through YOLOv12, in
164 identifying near-miss incidents on construction sites. The methodology encompasses frame-level video
165 sampling, model-based object detection, rule-based proximity analysis for near-miss identification, and
166 comparative evaluation using precision, recall, and spatial alignment metrics.

167 **Video Dataset and Frame Sampling:** The dataset comprised 38 publicly available construction-
168 related videos retrieved from the Pexels construction-video search page. The exact source-video set is
169 reported in Supplementary Table S1 using the downloaded Pexels video identifiers and filenames. The
170 source videos included HD and UHD footage with original frame rates ranging from 24 to 60 fps and
171 individual durations ranging from under 1 minute to approximately 2 minutes. Across the 38 videos,
172 the combined duration was approximately 903.13 seconds, or 15 minutes and 3 seconds, yielding
173 approximately 903 sampled frames at the 1 fps sampling rate. Because the source platform does not
174 provide standardized research metadata describing scene type, camera placement, lighting, weather,
175 occlusion, or worker-equipment interaction conditions, these characteristics were manually reviewed
176 by the authors through visual inspection of the videos. This review confirmed that the video set
177 included diverse construction and work-zone scenarios, including road resurfacing, coastal excavation,
178 urban equipment operation, and manual equipment handling, with visible variation in camera angle,
179 object scale, worker-equipment proximity, background clutter, motion blur, and partial occlusion. To
180 standardize evaluation across videos and reduce temporal redundancy, all videos were processed at a
181 normalized 1920×1080 resolution and sampled at one frame per second. Under this sampling strategy,
182 each second of video contributed one evaluation frame. Each extracted frame served as a consistent
183 input across all models.

184 **Model Selection and Proximity-Based Near-Miss Detection:** A total of 24 YOLO model variants,
185 spanning five generations from YOLOv8 to YOLOv12, were evaluated. Each generation included
186 configurations from lightweight nano (n) to extra-large (x). The Ultralytics interface handled model
187 loading, and all models were evaluated using the same inference pipeline across the complete set of
188 sampled frames. No model was fine-tuned on the study videos; all YOLO variants were evaluated
189 using their publicly available pretrained weights. Therefore, the study compares pretrained model
190 behavior under a controlled evaluation pipeline rather than comparing models retrained specifically for
191 this dataset. To ensure consistency, each YOLO variant was applied to the same sampled frames, with
192 the same normalized resolution, object classes, confidence settings, IoU matching criteria, centroid-
193 distance calculations, and near-miss matching rules. No model-specific threshold tuning, preprocessing
194 adjustment, or post-processing adjustment was applied. Therefore, differences in near-miss
195 performance reflect differences in model outputs under a controlled evaluation pipeline rather than
196 differences in preprocessing, post-processing, or evaluation criteria.

197 All models were evaluated under a shared inference configuration. Object detections were retained
198 using a fixed confidence threshold of 0.50, and non-maximum suppression was applied using an IoU
199 threshold of 0.45 across all YOLO variants before computing precision, recall, and F1 score. The
200 detector output was restricted to safety-relevant classes of interest: person, helmet, vest, truck,
201 excavator, crane, vehicle, machinery, and backhoe. Near-miss events were defined only between
202 person detections and equipment detections belonging to the subset {truck, excavator, crane, vehicle,
203 machinery, backhoe}. For each retained bounding box $b = (x_1, y_1, x_2, y_2)$, the centroid coordinates
204 were computed as:

205
$$c_x = \frac{x_1 + x_2}{2}$$

206
$$c_y = \frac{y_1 + y_2}{2}$$

207 The normalized centroid coordinates were computed as:

208
$$\hat{c}_x = \frac{c_x}{W}$$

209
$$\hat{c}_y = \frac{c_y}{H}$$

210 where $W = 1920$ and $H = 1080$ represent the standardized frame width and height, respectively.

211 The Euclidean image-space distance between a person centroid c_p and an equipment centroid c_e was
212 computed as:

213
$$d(p, e) = \sqrt{(c_{x,p} - c_{x,e})^2 + (c_{y,p} - c_{y,e})^2}$$

214 A person-equipment pair was classified as a near-miss candidate when $d(p, e)$ was less than or equal
215 to the selected proximity threshold. The primary analysis used a 250-pixel centroid-distance threshold,
216 and the same computation was repeated at 150, 200, and 300 pixels to assess threshold sensitivity.

217 A 250-pixel threshold was selected as the primary image-space proximity threshold for defining near-
218 miss events. Because the analyzed public construction footage did not include complete camera
219 calibration metadata, including focal length, mounting height, viewing angle, lens distortion, and
220 ground-plane reference measurements, the 250-pixel threshold should not be interpreted as a
221 universally calibrated physical distance. Instead, it functions as an operational threshold for comparing
222 model behavior under a fixed video resolution and camera perspective. For deployment in site-specific
223 monitoring systems, this threshold should be recalibrated using known reference distances, camera
224 calibration, homography-based ground-plane mapping, multi-camera geometry, or depth-aware
225 sensing. Because centroid distance is measured on the two-dimensional image plane, the same physical
226 worker-equipment separation can produce different pixel distances depending on camera height,
227 viewing angle, object depth, object size, lens distortion, and perspective projection. Therefore, the
228 proposed proximity rule should be interpreted as an image-space screening criterion for comparative
229 benchmarking rather than a direct physical measure of near-miss severity.

230 **Annotation Protocol and Performance Evaluation:** The framework generates CSV-based logs of
231 object detections and near-miss events for each model. These logs support computation of object
232 detection (OD) metrics, including precision, recall, and F1 score, using an IoU threshold of ≥ 0.3 . The
233 IoU threshold of 0.3 was selected because the evaluation videos contained small objects, partial
234 occlusion, motion blur, and variable camera perspectives, where stricter localization thresholds may
235 penalize detections that remain sufficient for downstream centroid-based interaction screening. This
236 threshold was applied uniformly across all models and was used for comparative object-detection
237 evaluation rather than for claiming fine-grained localization accuracy. The analysis compares model
238 outputs with a human-validated consensus ground truth. YOLOv8x detections provide preliminary
239 candidate annotations to reduce manual labeling effort and do not serve as final labels. A human-in-

240 the-loop annotation strategy mitigates circular validation bias associated with model-assisted pre-
241 labeling. Three domain experts, Ammar Alzarrad (Ph.D., M.ASCE, Construction Safety), Husnu S.
242 Narman (Ph.D., Computer Vision), and Sudipta Chowdhury (Ph.D., Industrial Safety), independently
243 review, correct, and supplement the annotations. The review process removes false-positive detections
244 caused by background clutter, corrects bounding-box localization errors, adds missing detections
245 absent from YOLOv8x outputs, and verifies person-equipment interaction labels used for near-miss
246 evaluation. Consensus discussion resolves disagreements in object classification, bounding-box
247 placement, and interaction boundaries until unanimous agreement is achieved for each contested frame.

248 Annotation guidelines were defined before expert review to ensure consistency across reviewers. A
249 “person” annotation included any visible worker or pedestrian within the sampled construction-related
250 frame. “Equipment” was operationally defined as visible construction or work-zone machinery and
251 vehicles capable of producing a worker-equipment proximity hazard, including trucks, excavators,
252 cranes, backhoes, and other mobile vehicles or machinery such as loaders, rollers, pavers, and forklifts
253 when visible. When equipment types were not represented as explicit YOLO class names, they were
254 assigned to the closest retained equipment superclass, such as “vehicle” or “machinery,” for
255 consistency with the model-output filtering rules. Static background objects, construction materials,
256 signs, cones, and distant or fully occluded machinery not involved in a potential worker-equipment
257 interaction were not labeled as near-miss equipment. A near-miss-positive interaction was defined as
258 a person-equipment pair whose bounding-box centroids, after review, satisfied the selected image-
259 space proximity rule. The final consensus benchmark contained approximately 903 sampled frames,
260 including 18 near-miss-positive sampled frames. These positive frames contained 67 consensus
261 person-equipment near-miss interaction instances under the 250-pixel image-space proximity rule. The
262 frame-level positive count was used for frame-level recall, while the interaction-instance count was
263 used for interaction-level near-miss precision, recall, and F1 evaluation. To reduce dependence on
264 YOLOv8x proposals, expert review was conducted in multiple passes: first, candidate YOLOv8x
265 detections were checked for false positives and localization errors; second, each full frame was visually
266 scanned for missed persons, equipment, and worker-equipment pairs absent from the candidate
267 annotations; third, centroid-based proximity labels were verified for eligible person-equipment pairs;
268 and fourth, disagreements were resolved through consensus discussion until a final human-corrected
269 reference label set was produced.

270 To quantify annotation reliability before consensus resolution, inter-rater agreement was computed
271 using a stratified random sample of 200 frames, selected to represent different video conditions and
272 interaction types. Annotators independently labeled frame-level near-miss presence according to the
273 predefined proximity rule. Agreement among the three annotators was measured using Fleiss’ kappa,
274 yielding $\kappa = 0.87$, indicating strong agreement for a safety-critical visual labeling task. Pairwise
275 agreement rates exceeded 90% across all annotator combinations. The finalized benchmark was
276 therefore a human-corrected, consensus-resolved ground truth dataset rather than the raw YOLOv8x
277 output. This consensus dataset served as the reference for all comparative evaluations, including
278 YOLOv8x itself. Therefore, the analysis does not compare YOLOv8x against its raw predictions but
279 evaluates it against a revised ground truth that removes false positives, corrects localization errors, and
280 adds missing expert-identified detections. This procedure reduces circular validation bias and ensures
281 that performance metrics reflect model behavior against an independently reviewed benchmark rather
282 than artifacts of the initial pre-labeling model.

283 Near-miss (NM) performance was evaluated using centroid normalization and fuzzy matching to assess
284 detection accuracy at both the instance and frame levels. Fuzzy matching between predicted and
285 reference near-miss events was performed within each frame using greedy nearest-neighbor one-to-

286 one matching. A predicted person-equipment interaction was matched to a reference interaction when
287 the corresponding centroid distance was within a 50-pixel tolerance, with ties resolved by selecting the
288 smallest centroid distance first. Each reference interaction could be matched to at most one predicted
289 interaction. Additional metrics included frame-level recall, average IoU within expanded person-
290 equipment interaction zones, and mean proximity distance across all interactions. For the interaction-
291 zone IoU calculation, each person and equipment bounding box was expanded outward by 50 pixels
292 on all sides. The intersection-over-union was then computed between the expanded person box and the
293 expanded equipment box using the standard intersection area divided by union area formula. If the
294 union area was zero, the IoU was assigned a value of 0.0. Visualization of results used bar plots, line
295 charts, and confusion matrices. The final evaluation also included an internal consistency check to
296 verify that raw near-miss event logs aligned with the computed NM metrics. To clarify the unit of
297 analysis for the statistical comparison, the ANOVA was not conducted using a single aggregate NM
298 F1 score per model. Instead, a per-frame NM F1 score was computed for each of the 18 near-miss-
299 positive sampled frames by comparing model-predicted person-equipment interactions with the
300 consensus interactions in that frame using the fuzzy matching procedure described above. Each of the
301 24 YOLO variants therefore contributed 18 frame-level NM F1 records, yielding 432 model-frame
302 records. Model variant was treated as the grouping factor in a one-way ANOVA, followed by Tukey-
303 adjusted post hoc comparisons. Because these observations were sampled from video data and may
304 retain temporal dependence within source videos, the ANOVA was interpreted as an exploratory
305 within-benchmark comparison rather than as population-level statistical validation. Temporal
306 clustering by source video and repeated evaluation of the same frames across models were not modeled
307 as formal random effects. Therefore, no population-level inference about real construction-site
308 performance is claimed from the ANOVA results. The primary interpretation, therefore, emphasizes
309 descriptive model rankings, score differences, threshold-sensitivity patterns, and consistency checks.

310 **4 Results**

311 The evaluation framework computes object detection (OD) and near-miss (NM) performance metrics
312 for each model. Within this benchmark, a near-miss candidate was operationally defined when the
313 pixel distance between a detected person and a piece of equipment was less than or equal to 250 pixels.
314 In the final consensus benchmark, 18 near-miss-positive sampled frames were identified, containing
315 67 consensus person-equipment near-miss interaction instances under the 250-pixel image-space
316 proximity rule. Figure 1 illustrates sample near-miss detections from construction site footage,
317 highlighting hazardous person-equipment interactions across diverse scenarios, including road
318 resurfacing, coastal excavation, urban operations, and manual equipment handling. These examples
319 capture varied public-video conditions, including differences in resolution, lighting, motion blur, and
320 camera angles, illustrating the pipeline’s comparative behavior under visually diverse image-space
321 conditions rather than validating robustness across all real jobsite deployments.

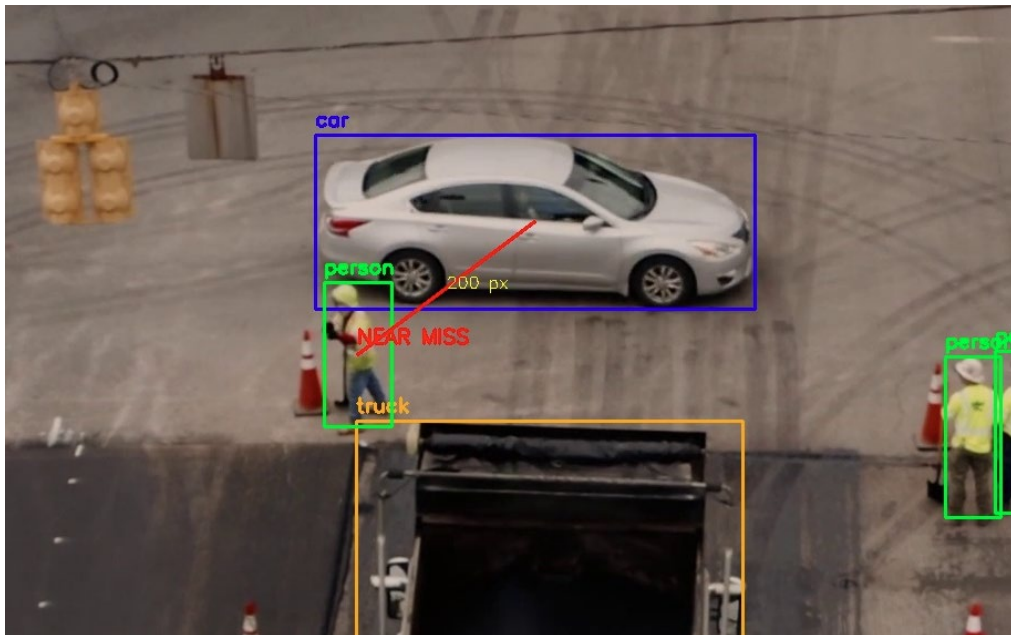
322 YOLOv8x achieved the highest object-detection F1 score (0.980) and the highest near-miss F1 score
323 (0.930). As shown in Table 2, YOLOv8x obtained NM precision = 0.819, NM recall = 0.879, frame-
324 level recall = 1.000, average person-equipment interaction-zone IoU = 0.256, and detected near-miss-
325 positive frames = 18. Other strong performers included YOLOv8l (OD F1 = 0.953, NM F1 = 0.879),
326 YOLOv9c (NM F1 = 0.863), and YOLOv11x (NM F1 = 0.860). YOLOv10b and YOLOv12l also
327 performed above the field median, each achieving an NM F1 score of 0.820. As shown in Figure 2,
328 YOLOv8x and YOLOv8l led the near-miss F1 rankings, followed by selected YOLOv9 and YOLOv11
329 variants. The declining trend beyond the top models illustrates the challenge of maintaining detection
330 consistency and spatial reasoning in compact architectures. Table 2 summarizes object-detection and
331 near-miss metrics for all 24 YOLO variants, ranked by near-miss F1 score.

Table 2. Object detection and near-miss performance across 24 YOLO models

Model	Version	OD F1	NM Precision	NM Recall	NM F1	Frame-Level Recall	Avg IOU (Person–Equip Zone)	Detected NM Frames
yolov8x	YOLOv8	0.980	0.819	0.879	0.930	1.000	0.256	18
yolov8l	YOLOv8	0.953	0.858	0.848	0.879	0.944	0.280	17
yolov9c	YOLOv9	0.923	0.856	0.813	0.863	0.889	0.281	18
yolov11x	YOLOv11	0.927	0.830	0.874	0.860	0.833	0.281	16
yolov8s	YOLOv8	0.804	0.817	0.835	0.831	0.833	0.278	17
yolov8n	YOLOv8	0.741	0.827	0.820	0.831	0.056	0.173	11
yolov11s	YOLOv11	0.832	0.826	0.829	0.828	0.778	0.253	15
yolov9m	YOLOv9	0.929	0.819	0.823	0.822	0.889	0.268	17
yolov12l	YOLOv12	0.911	0.827	0.822	0.820	0.889	0.237	18
yolov9t	YOLOv9	0.720	0.826	0.814	0.820	0.167	0.185	9
yolov10n	YOLOv10	0.713	0.819	0.812	0.820	0.278	0.171	13
yolov10b	YOLOv10	0.920	0.814	0.803	0.820	0.944	0.269	18
yolov9e	YOLOv9	0.936	0.840	0.817	0.819	1.000	0.229	18
yolov12x	YOLOv12	0.942	0.803	0.824	0.817	1.000	0.252	18
yolov10m	YOLOv10	0.899	0.819	0.826	0.816	0.833	0.248	16
yolov10l	YOLOv10	0.916	0.815	0.823	0.816	0.944	0.252	18
yolov8m	YOLOv8	0.910	0.855	0.839	0.815	0.944	0.266	18
yolov11n	YOLOv11	0.733	0.817	0.818	0.815	0.056	0.159	7

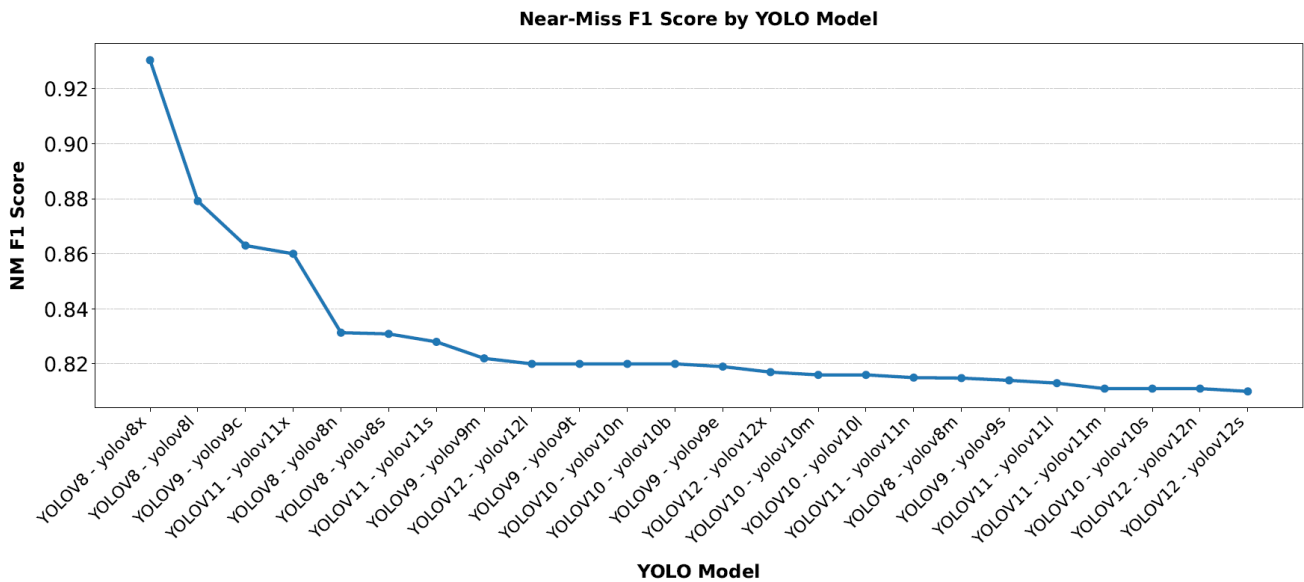
yolov9s	YOLOv9	0.836	0.814	0.814	0.814	0.889	0.296	18
yolov11l	YOLOv11	0.900	0.816	0.821	0.813	0.833	0.217	16
yolov11m	YOLOv11	0.898	0.812	0.814	0.811	1.000	0.247	18
yolov10s	YOLOv10	0.843	0.822	0.813	0.811	0.944	0.276	18
yolov12n	YOLOv12	0.754	0.825	0.825	0.811	0.222	0.188	12
yolov12s	YOLOv12	0.836	0.825	0.814	0.810	0.944	0.293	18

333 OD F1 is the object detection F1 score against the human-validated consensus ground truth initialized
 334 from YOLOv8x candidate annotations. NM Precision, NM Recall, and NM F1 are computed for near-
 335 miss events under the 250-pixel proximity rule. Frame-level recall measures the proportion of frames
 336 with at least one correctly detected near-miss, Avg IOU is the mean intersection-over-union within
 337 person-equipment interaction zones, and Detected NM Frames is the number of sampled frames in
 338 which the model flagged at least one near-miss candidate. This frame-level count is distinct from the
 339 67 consensus person-equipment near-miss interaction instances used for interaction-level precision,
 340 recall, and F1 evaluation.



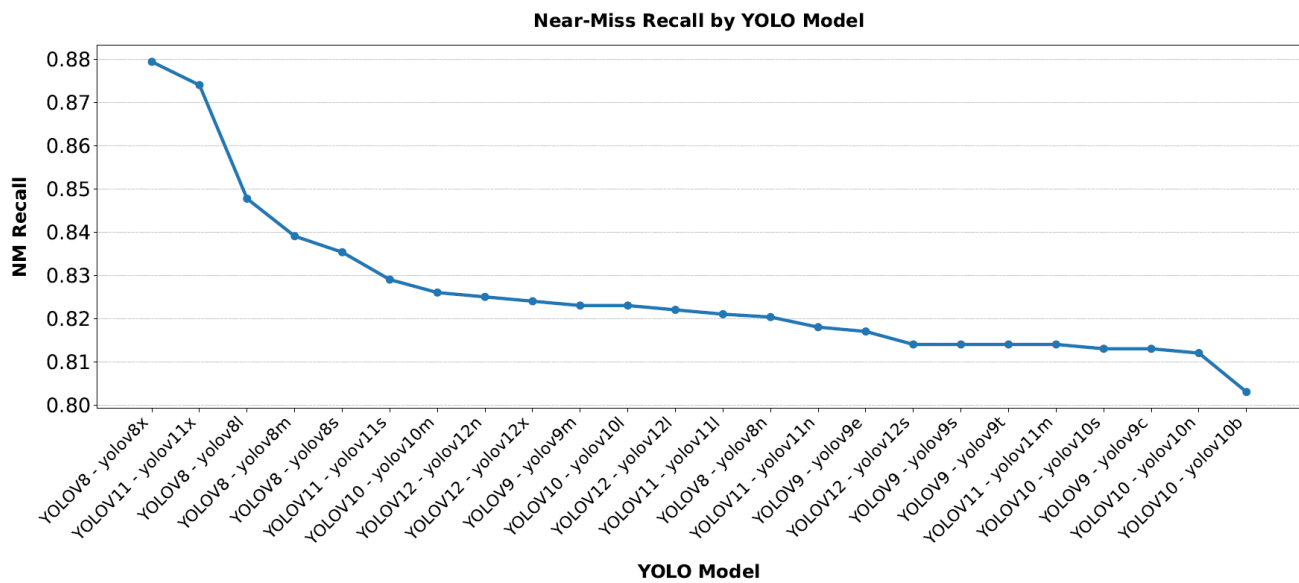


341 **Figure 1: Annotated frames depicting detected near-miss interactions**



342
343 **Figure 2: Near-Miss F1 Score by YOLO Model**

344 Figure 3 illustrates differences in near-miss recall across YOLO variants. The leading models
 345 maintained stronger near-miss recall and more stable frame-level object coverage, suggesting that
 346 detection consistency is important for near-miss identification under cluttered and partially occluded
 347 construction-site conditions.

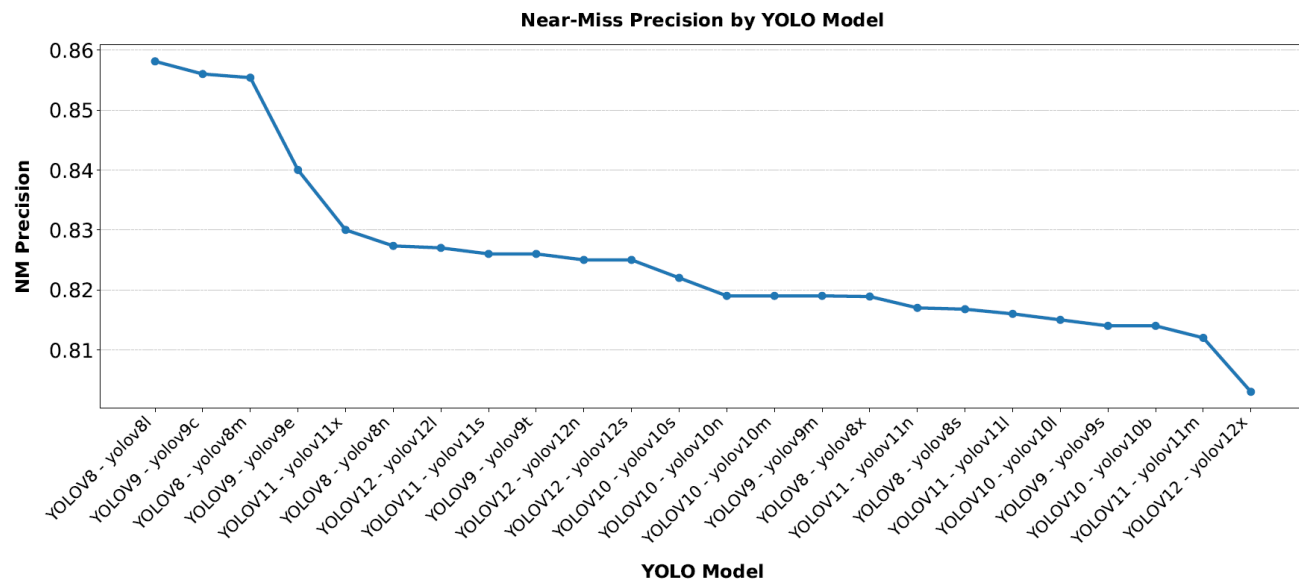


348

349

Figure 3: Near-Miss Recall by YOLO Model

350 As shown in Figure 4, near-miss precision was comparatively stable across several models, indicating
 351 that many variants produced spatially plausible near-miss detections when relevant objects were
 352 successfully detected. However, precision should be interpreted alongside frame-level recall because
 353 models with limited frame coverage may appear precise on the smaller subset of frames where
 354 detections occur, while still missing important near-miss events.



355

356

Figure 4: Near-Miss Precision by YOLO Model

357 A one-way ANOVA was conducted on observation-level NM F1 records rather than on a single
 358 aggregate NM F1 score per model. Each of the 24 YOLO variants contributed 18 near-miss-positive
 359 frame-level observations, yielding 432 model-observation records. The analysis showed significant
 360 variation in NM F1 across model variants, $F(23, 408) = 416.62, p < 0.001$. Tukey-adjusted post hoc

361 comparisons indicated that YOLOv8x remained separated from the lower-performing model cluster,
362 defined as models with aggregate NM F1 scores at or below 0.820. These exploratory statistical results
363 are consistent with the descriptive benchmark ranking, where YOLOv8x achieved the highest
364 aggregate NM F1 score, followed by YOLOv8l, YOLOv9c, and YOLOv11x. However, because the
365 observations were sampled from video footage and may retain temporal dependence within source
366 videos, the ANOVA should be interpreted as exploratory evidence within the sampled benchmark
367 rather than broad population-level statistical validation. Temporal clustering by source video and
368 repeated evaluation of the same positive frames across models were not modeled as random effects,
369 and the 18 positive frames represent a limited positive-event sample; therefore, the ANOVA results
370 should not be interpreted as population-level evidence of real-world construction-site performance.
371 Although YOLOv8x was used to generate the initial candidate annotations, its final performance was
372 evaluated against the same human-corrected consensus ground truth used for all other models. This
373 distinction is important because the final benchmark incorporated expert-added missed detections,
374 corrected bounding boxes, and removed false positives rather than relying on unmodified YOLOv8x
375 outputs.

376 Several lower-performing models, including YOLOv12s (NM F1 = 0.810), YOLOv12n (NM F1 =
377 0.811), YOLOv10s (NM F1 = 0.811), and YOLOv11m (NM F1 = 0.811), trailed the top-performing
378 model by more than 0.10 NM F1 points. Models with particularly low frame-level recall, such as
379 YOLOv8n (Frame-Level Recall = 0.056, NM F1 = 0.831) and YOLOv11n (Frame-Level Recall =
380 0.056, NM F1 = 0.815), correctly identified near-miss-positive frames less consistently. However,
381 when these models produced detections, their proximity-based classifications remained broadly
382 comparable to those of larger variants. YOLOv9t showed a similar pattern, with Frame-Level Recall
383 = 0.167, NM F1 = 0.820, and only 9 detected near-miss-positive frames. This suggests that its lower
384 near-miss coverage was primarily due to limited frame-level object detection rather than a failure of
385 the proximity logic itself. Its object-detection F1 score was also relatively low at 0.720.

386 Correlation analysis showed a moderate positive relationship between NM F1 and frame-level recall
387 ($r = 0.69$), compared with a weaker relationship between NM F1 and detected near-miss-positive frame
388 count ($r = 0.35$). The correlation between NM F1 and average proximity distance was negligible ($r = -$
389 0.04), suggesting that detection coverage and recall were more influential for near-miss performance
390 than proximity distance alone. Internal consistency checks were conducted to verify that the frame-
391 level detected near-miss counts and interaction-level NM precision, recall, and F1 scores were
392 computed from the corresponding model event logs. Across all 24 models, OD F1 scores ranged from
393 0.713 to 0.980, NM F1 scores ranged from 0.810 to 0.930, frame-level recall ranged from 0.056 to
394 1.000, and average person-equipment interaction-zone IoU ranged from 0.159 to 0.296. As
395 summarized in Table 3, model rankings remained broadly stable across proximity thresholds from 150
396 to 300 pixels, indicating that the comparative ranking was not driven solely by the selected 250-pixel
397 cutoff. However, this threshold-sensitivity result does not remove the inherent limitations of centroid-
398 based image-space proximity, because pixel distances remain affected by perspective, camera
399 geometry, object depth, and scale variation. Accordingly, the reported results should be interpreted as
400 comparative model performance rather than absolute physical distance measurements.

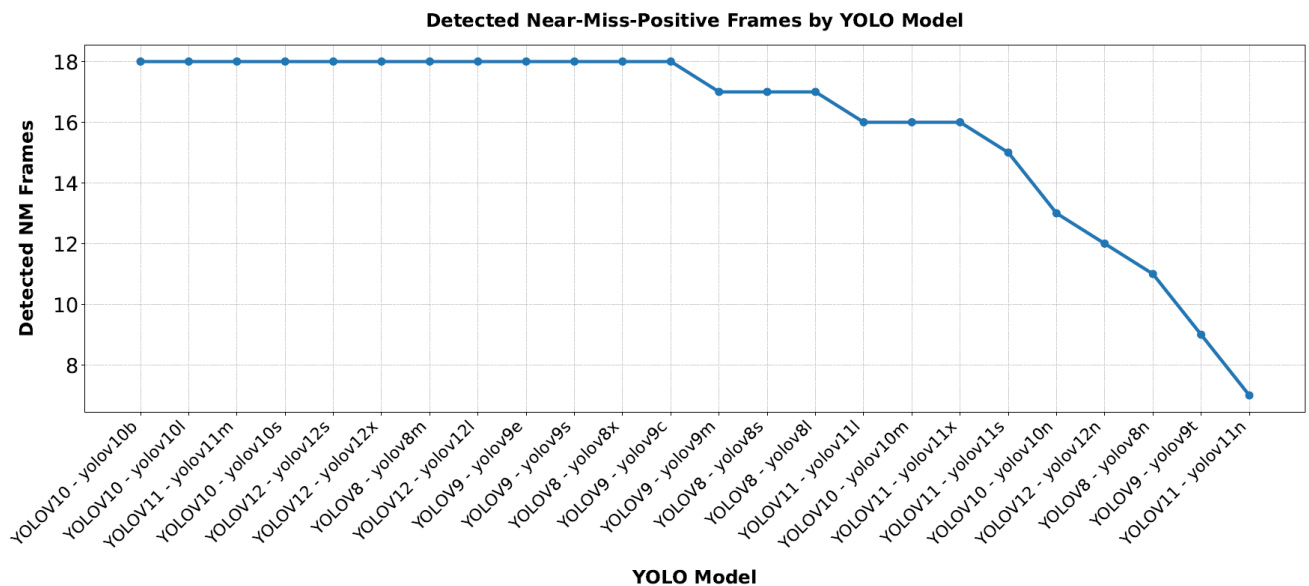
401 To assess sensitivity to the selected proximity threshold, the analysis recalculates near-miss F1 scores
402 at 150, 200, 250, and 300 pixels for all 24 YOLO variants. Table 3 presents the results. The relative
403 model-ranking pattern remains broadly stable across thresholds, indicating that the main findings do
404 not depend solely on the 250-pixel cutoff.

405 **Table 3. Threshold-Sensitivity Analysis of Near-miss F1 Scores Across Proximity Thresholds**

Model	NM F1 @150px	NM F1 @200px	NM F1 @250px	NM F1 @300px
YOLOv8x	0.901	0.918	0.930	0.924
YOLOv8l	0.850	0.866	0.879	0.872
YOLOv9c	0.835	0.850	0.863	0.856
YOLOv11x	0.832	0.847	0.860	0.853
YOLOv8n	0.803	0.818	0.831	0.824
YOLOv8s	0.802	0.817	0.831	0.823
YOLOv11s	0.800	0.815	0.828	0.822
YOLOv9m	0.794	0.809	0.822	0.815
YOLOv12l	0.792	0.806	0.820	0.813
YOLOv9t	0.789	0.805	0.820	0.812
YOLOv10n	0.787	0.804	0.820	0.811
YOLOv10b	0.790	0.805	0.820	0.812
YOLOv9e	0.791	0.805	0.819	0.813
YOLOv12x	0.789	0.803	0.817	0.810
YOLOv10m	0.788	0.802	0.816	0.810
YOLOv10l	0.787	0.801	0.816	0.809
YOLOv11n	0.786	0.800	0.815	0.808
YOLOv8m	0.787	0.801	0.815	0.809
YOLOv9s	0.784	0.799	0.814	0.807
YOLOv11l	0.783	0.798	0.813	0.807
YOLOv11m	0.781	0.796	0.811	0.805
YOLOv10s	0.782	0.797	0.811	0.806

YOLOv12n	0.779	0.795	0.811	0.803
YOLOv12s	0.778	0.794	0.810	0.802

406 Figure 5 displays the number of sampled frames in which each YOLO model flagged at least one near-
407 miss candidate. This frame-level count reflects coverage across the 18 consensus near-miss-positive
408 sampled frames, but it is distinct from the 67 consensus person-equipment interaction instances used
409 for interaction-level precision, recall, and F1 evaluation. Models with higher detected-frame counts
410 generally provided broader frame-level coverage, while models with lower detected-frame counts may
411 have missed near-miss-positive frames. Therefore, detected-frame count should be interpreted
412 alongside precision, recall, F1 score, and frame-level recall rather than as a standalone indicator of
413 model effectiveness.



414
415 **Figure 5: Detected Near-Miss-Positive Frames by Model Across YOLO Versions**

416 **5 Discussion**

417 This study evaluated 24 YOLO model variants, spanning YOLOv8 through YOLOv12, for interaction-
418 level near-miss detection in construction environments. The results show that stronger object-detection
419 consistency and stable frame-level coverage are important for identifying hazardous worker-equipment
420 interactions. Among the evaluated models, YOLOv8x achieved the highest near-miss performance
421 (NM F1 = 0.930), followed by YOLOv8l (NM F1 = 0.879) and YOLOv9c (NM F1 = 0.863). These
422 results indicate that larger or higher-capacity models generally performed better than lightweight
423 variants in this task, particularly under visually complex conditions involving clutter, partial occlusion,
424 motion blur, and varying object scales. These findings align with prior construction safety studies that
425 have demonstrated the effectiveness of YOLO-based and deep-learning-based approaches for tasks
426 such as drowsiness detection, PPE compliance, unsafe-behavior recognition, and safety-compliance
427 monitoring (Onososen et al., 2025; Mohy et al., 2025; Alotaibi and Ma, 2025). However, the present
428 study extends this body of work by focusing on interaction-level near-miss detection rather than
429 isolated object or behavior recognition. In near-miss detection, accurate object localization alone is not
430 sufficient; the model must also produce stable detections that support reliable spatial reasoning between

431 workers and equipment. This distinction is important because a missed person or equipment detection
432 can directly affect whether a hazardous interaction is identified.

433 The performance of lightweight models further highlights the importance of frame-level coverage. For
434 example, YOLOv8n and YOLOv11n both produced very low frame-level recall values of 0.056, while
435 YOLOv9t achieved a frame-level recall of 0.167. These models flagged fewer near-miss events than
436 higher-recall variants, suggesting that their limitations were primarily associated with inconsistent
437 detection coverage rather than failure of the centroid-based proximity logic itself. This pattern is
438 consistent with prior research showing that occlusion, background clutter, and complex surveillance
439 conditions can reduce the reliability of object detection in real-world visual monitoring systems
440 (Ouardirhi et al., 2024). For construction safety applications, these results suggest that model capacity
441 should be considered not only in terms of object-detection accuracy but also in terms of whether the
442 model can maintain reliable detections across frames. Correlation analysis provided additional insight
443 into the factors influencing near-miss performance. Frame-level recall showed a moderate positive
444 relationship with NM F1 performance ($r = 0.69$), while the relationship between NM F1 and average
445 proximity distance was negligible ($r = -0.04$). This suggests that detection consistency and frame
446 coverage were more influential for near-miss identification than the absolute closeness of detected
447 objects alone. The threshold-sensitivity analysis also showed that model rankings remained broadly
448 stable across proximity thresholds from 150 to 300 pixels. Therefore, the main comparative findings
449 were not driven solely by the selected 250-pixel threshold. At the same time, because the proximity
450 threshold was defined in image space rather than through full camera calibration, the results should be
451 interpreted as comparative model performance rather than absolute physical distance measurements.
452 Future deployment should incorporate camera calibration, ground-plane mapping, homography-based
453 distance estimation, or depth-aware sensing to translate image-space proximity into physically
454 meaningful worker-equipment distances. This limitation is especially important for construction scenes
455 because workers and equipment may appear close in the image while being separated in depth, or may
456 appear farther apart in pixels despite being physically close on the ground plane. As a result, centroid-
457 based image-space proximity should not be interpreted as a complete near-miss definition or as a
458 substitute for calibrated physical-distance estimation. Its role in this study is to provide a consistent
459 operational rule for comparing detector behavior under the same image-space conditions.

460 The results also show that newer YOLO versions did not automatically outperform earlier variants.
461 Although YOLOv10, YOLOv11, and YOLOv12 include architectural updates, their performance
462 varied across model sizes. For example, YOLOv12l performed competitively, but YOLOv12n and
463 YOLOv12s were among the lower-performing variants in near-miss F1. Similarly, YOLOv9t showed
464 limited frame-level coverage despite belonging to a newer model generation. These results suggest that
465 improvements in general object-detection architectures do not always translate directly into better
466 performance for spatial proximity-based near-miss detection. Task-specific evaluation is therefore
467 necessary when selecting models for construction safety monitoring, especially when the downstream
468 goal involves interaction-level risk assessment rather than object detection alone. For practitioners, the
469 findings suggest that model selection should depend on deployment context rather than accuracy alone.
470 Within the present image-space benchmark, mid-to-large models such as YOLOv8x, YOLOv8l, and
471 YOLOv9c achieved the strongest near-miss performance and may be preferable when the priority is
472 reducing missed hazardous interactions. However, these models typically require greater
473 computational resources, memory, and inference time than lightweight variants, which may limit their
474 suitability for low-power edge devices, embedded cameras, or real-time on-site deployment. Smaller
475 models may be more practical for resource-constrained environments, but their lower frame-level
476 coverage may increase the risk of missed near-miss events. Therefore, deployment decisions should
477 consider hardware availability, acceptable latency, camera density, network bandwidth, and whether

478 inference is performed locally at the edge or centrally on a server. However, this study did not directly
479 measure inference latency, memory consumption, energy use, hardware-specific throughput, or real-
480 time edge-device performance. Therefore, the deployment discussion should be interpreted as a
481 practical consideration based on model size and detection behavior, not as a hardware-benchmarking
482 result.

483 The tradeoff between false positives and false negatives is also central to practical use. False negatives
484 are especially concerning in construction safety because they represent missed hazardous interactions
485 and may prevent timely intervention. False positives, while generally less severe, can contribute to
486 alarm fatigue, reduce worker trust, and increase the review burden for safety personnel. The results
487 show that detected near-miss-positive frame count alone is not a sufficient indicator of model
488 effectiveness. Models that flag more near-miss-positive frames may provide broader frame-level
489 coverage but may also introduce additional false positives at the interaction level, while models with
490 low frame-level recall may miss important hazardous frames. Accordingly, practical systems should
491 tune detection and alert thresholds according to site-specific risk tolerance. High-risk zones involving
492 heavy equipment may prioritize recall to reduce missed events, whereas lower-risk monitoring contexts
493 may require stricter precision to avoid excessive alerts.

494 The external validity of the benchmark is limited by the use of publicly available Pexels videos rather
495 than purpose-collected operational jobsite footage. Although the 38 videos varied in camera angle,
496 object scale, worker-equipment proximity, clutter, motion blur, and partial occlusion, the dataset does
497 not fully capture the frequency, severity, camera geometry, equipment movement patterns, weather
498 conditions, site-specific work practices, or hazard distributions encountered in active construction
499 monitoring. The absence of standardized metadata for camera placement, weather, lighting, and site
500 operations also limits the extent to which the results can be generalized beyond the sampled video set.
501 Therefore, the findings should be interpreted as comparative model behavior under selected image-
502 space construction-video conditions, not as a broadly validated near-miss detection solution for field
503 deployment. In addition, because the benchmark contained only 18 near-miss-positive sampled frames
504 and 67 consensus near-miss interaction instances, the near-miss performance estimates may be
505 sensitive to a small number of missed or additional detections and should be interpreted as exploratory
506 rather than definitive. Future work should validate the framework using longer, purpose-collected, site-
507 specific videos with documented camera geometry, calibrated distances, operational metadata, and
508 safety-expert-labeled near-miss events.

509 Overall, this study demonstrates the value of evaluating construction safety models beyond
510 conventional object-detection metrics. Near-miss detection requires stable object recognition, spatial
511 reasoning, threshold sensitivity assessment, and careful interpretation of deployment constraints. The
512 proposed framework provides a reproducible basis for comparing YOLO model variants in interaction-
513 level safety monitoring. Future work should expand the dataset with longer and more diverse site-
514 specific videos, incorporate calibrated physical-distance estimation, directly benchmark real-time
515 latency, memory use, and throughput on edge and server hardware, and develop cost-sensitive metrics
516 that separately weight false positives and false negatives according to construction safety priorities.

517 **6 Conclusion**

518 This research developed and applied an evaluation framework to assess 24 pretrained YOLO models,
519 from YOLOv8 to YOLOv12, for construction near-miss detection. By combining object detection
520 outputs with centroid-based proximity analysis, the framework examined how well different YOLO
521 variants could identify potentially hazardous worker-equipment interactions in publicly available

522 construction-related video footage. The results indicate that near-miss detection depends heavily on
523 consistent frame-level object coverage. YOLOv8x, YOLOv8l, and YOLOv9c produced the strongest
524 overall performance, while several lightweight models showed reduced reliability in detecting near-
525 miss events across frames. This suggests that model size and detection stability are important
526 considerations when applying computer vision systems to interaction-level safety monitoring. The
527 study also demonstrates the need to evaluate safety-monitoring models beyond standard object-
528 detection metrics. Near-miss detection requires not only accurate object recognition, but also stable
529 localization and reliable spatial interpretation of worker-equipment relationships. The proposed
530 framework offers a reproducible basis for comparing YOLO-based safety systems. It can support future
531 work on calibrated distance estimation, temporal tracking, direct latency and edge-device
532 benchmarking, and cost-sensitive evaluation of false positives and false negatives.

533 **7 Conflict of Interest**

534 The authors declare that the research was conducted in the absence of any commercial or financial
535 relationships that could be construed as a potential conflict of interest.

536 **8 Author Contributions**

537 SA: Conceptualization, Visualization, Data curation, Methodology, Validation, Writing-original, draft,
538 Writing-review and editing. AA: Methodology, Funding acquisition, Project administration, Writing-
539 original draft. SC: Methodology, Software, Validation, Visualization, Writing-original draft. HN:
540 Formal Analysis, Investigation, Writing-original draft.

541 **9 Funding**

542 The authors declare that financial support was not received for the research, authorship, and/or
543 publication of this article.

544 **10 References**

- 545 Ait Abdelali, H., Derrouz, H., Zennayi, Y., Thami, R. O. H., and Bourzeix, F. 2021. "Multiple
546 hypothesis detection and tracking using deep learning for video traffic surveillance." *IEEE Access*
547 9, 164282-164291. <https://doi.org/10.1109/ACCESS.2021.3133529>.
- 548 Alotaibi, R. T. T., and Ma, S. 2025. "Real-time detection of unsafe worker behaviors via adaptive
549 vision transformers in construction sites." *Buildings* 15, 4205.
550 <https://doi.org/10.3390/buildings15224205>.
- 551 Alzarrad, A., Emanuels, C., Imtiaz, M., and Akbar, H. 2021. "Automatic assessment of buildings'
552 location fitness for solar panel installation using drones and neural networks." *CivilEng* 2, 1052-
553 1064. <https://doi.org/10.3390/civileng2040056>.
- 554 Alzarrad, A., Awolusi, I., Hatamleh, M. T., and Terreno, S. 2022. "Automatic assessment of roof
555 conditions using artificial intelligence and unmanned aerial vehicles." *Front. Built Environ.* 8,
556 1026225. <https://doi.org/10.3389/fbuil.2022.1026225>.
- 557 Alzarrad, A., Ali, S., Chowdhury, S., and Narman, H. S. 2025. "Reducing gun-related incidents on
558 construction sites: An AI-driven approach for automated detection of people, tools, and firearms."
559 In: *Computing in Civil Engineering 2025*. <https://doi.org/10.1061/9780784486436.032>.

- 560 Damai, A., Song, H., Narman, H. S., Lambert, A., and Alzarrad, A. 2025. "Enhancing railway safety:
561 A machine learning approach for automated detection of missing track bolts." In: *Computing in*
562 *Civil Engineering 2025*. <https://doi.org/10.1061/9780784486436.024>.
- 563 Dzung, R.-J., Fan, B., and Hsieh, T.-L. 2025. "Dynamic collision alert system for collaboration of
564 construction equipment and workers." *Buildings* 15, 110.
565 <https://doi.org/10.3390/buildings15010110>.
- 566 Feng, R., Miao, Y., and Zheng, J. 2024. "A YOLO-based intelligent detection algorithm for risk
567 assessment of construction sites." *J. Intell. Constr.* 2, 1-18.
568 <https://doi.org/10.26599/JIC.2024.9180037>.
- 569 Guo, Z., Wang, C., Yang, G., Huang, Z., and Li, G. 2022. "MSFT-YOLO: Improved YOLOv5 based
570 on transformer for detecting defects of steel surface." *Sensors* 22, 3467.
571 <https://doi.org/10.3390/s22093467>.
- 572 Jiang, B., Zhu, Y., Zhou, Z., Zhang, Y., Xu, L., and Wang, J. 2024. "An improved safety belt
573 detection algorithm for high-altitude work based on YOLOv8." *Electronics* 13, 850.
574 <https://doi.org/10.3390/electronics13050850>.
- 575 Kizilay, F., Narman, M. R., Song, H., Narman, H. S., Cosgun, C., and Alzarrad, A. 2024. "Evaluating
576 fine-tuned deep learning models for real-time earthquake damage assessment with drone-based
577 images." *AI Civ. Eng.* 3, 15. <https://doi.org/10.1007/s43503-024-00034-6>.
- 578 Le, V. T., Song, H., Narman, H. S., Zhu, P., Alzarrad, A., Cisco, A., et al. 2025. "Automated
579 measurement of horizontal gauge deviation in railroads using depth sensor camera and machine
580 learning." *IEEE Access.* 13, pp. 215324-215338. <https://doi.org/10.1109/ACCESS.2025.3641797>.
- 581 Lee, Y.-R., Jung, S.-H., Kang, K.-S., Ryu, K.-C., and Ryu, H.-G. 2023. "Deep learning-based
582 framework for monitoring wearing personal protective equipment on construction sites." *J.*
583 *Comput. Des. Eng.* 10, 905-917. <https://doi.org/10.1093/jcde/qwad019>.
- 584 Lim, L. M., Yang, L., Zhu, W., Mohamed, A. S. A., and Ali, M. K. M. 2025. "Near miss detection
585 using distancing monitoring and distance-based proximal indicators." *IEEE Access* 13, 48449-
586 48468. <https://doi.org/10.1109/ACCESS.2025.3548108>.
- 587 López, L., Suárez-Ramírez, J., Alemán-Flores, M., and Monzón, N. 2025. "Automated PPE
588 compliance monitoring in industrial environments using deep learning-based detection and pose
589 estimation." *Autom. Constr.* 176, 106231. <https://doi.org/10.1016/j.autcon.2025.106231>.
- 590 Lyu, Y., Zhang, T., Li, X., Liu, A., and Shi, G. 2025. "LightUAV-YOLO: A lightweight object
591 detection model for unmanned aerial vehicle image." *J. Supercomput.* 81, 105.
592 <https://doi.org/10.1007/s11227-024-06611-x>.
- 593 Maulana, R. A., and Hardiansyah. 2025. "Implementation of YOLO algorithm for drowsiness
594 detection as an additional safety feature in crane operation." *J. Inotera* 10, 113-120.
595 <https://doi.org/10.31572/inotera.Vol10.Iss1.2025.ID458>.

- 596 Mohy, A. A., Bassioni, H. A., Elgendi, E. O., and Hassan, T. M. 2025. "Deep learning-enabled
597 computer vision model for automated safety compliance in construction environments." *J. Inf.*
598 *Technol. Constr.* 30, 1398-1430. <https://doi.org/10.36680/j.itcon.2025.057>.
- 599 Olimat, H., Alwashah, Z., Abudayyeh, O., and Liu, H. 2025. "Data-driven analysis of construction
600 safety dynamics: Regulatory frameworks, evolutionary patterns, and technological innovations."
601 *Buildings* 15, 1680. <https://doi.org/10.3390/buildings15101680>.
- 602 Onososen, A. O., Musonda, I., Onatayo, D., Saka, A. B., Adekunle, S. A., and Onatayo, E. 2025.
603 "Drowsiness detection of construction workers using YOLOv8 and computer vision techniques."
604 *Buildings* 15, 500. <https://doi.org/10.3390/buildings15030500>.
- 605 Ouairihi, Z., Mahmoudi, S. A., and Zbakh, M. 2024. "Enhancing object detection in smart video
606 surveillance: A survey of occlusion-handling approaches." *Electronics* 13, 541.
607 <https://doi.org/10.3390/electronics13030541>.
- 608 Saeedizadeh, N., Jalali, S. M. J., Khan, B., and Mohamed, S. 2025. "Cutting-edge deep learning
609 methods for image-based object detection in autonomous driving: In-depth survey." *Expert Syst.*
610 42, e70020. <https://doi.org/10.1111/exsy.70020>.
- 611 U.S. Bureau of Labor Statistics (BLS). 2026. "National Census of Fatal Occupational Injuries in
612 2024." News release, February 19. <https://www.bls.gov/news.release/pdf/foi.pdf>.
- 613 Yongxiang, G., Xingbin, L., and Xiaolin, Q. 2022. "YouTube-GDD: A challenging gun detection
614 dataset with rich contextual information." *arXiv*. <https://doi.org/10.48550/arXiv.2203.04129>.
- 615 Zhou, S., Zhou, H., and Qian, L. 2025. "A multi-scale small object detection algorithm SMA-YOLO
616 for UAV remote sensing images." *Sci. Rep.* 15, 9255. [https://doi.org/10.1038/s41598-025-92344-](https://doi.org/10.1038/s41598-025-92344-7)
617 [7](https://doi.org/10.1038/s41598-025-92344-7).
- 618 Zou, X., and Hu, Y. 2024. "Hidden danger detection and identification system of power transmission
619 tower based on YOLOv11." *Acad. J. Sci. Technol.* 13, 224-231. <https://doi.org/10.54097/rs28p954>.

620 **11 Data Availability Statement**

621 The video data used in this study were obtained from publicly available construction-related videos on
622 Pexels through the construction-video search page:
623 <https://www.pexels.com/search/videos/construction/>. Because the Pexels search results page is
624 dynamically updated and may be re-ranked over time, the exact benchmark source set is defined by the
625 38 Pexels video IDs and filenames listed in Supplementary Table S1. The frame-sampling procedure,
626 benchmark composition, annotation protocol, consensus-review process, and shared
627 inference/evaluation settings are described in the manuscript to support reproducibility. The authors
628 do not redistribute the underlying videos; readers should access them through Pexels, subject to the
629 platform's licensing and availability conditions. Additional implementation files or clarifications can
630 be requested from the corresponding author.