# M&M: Tackling False Positives in Mammography with a Multi-view and Multi-instance Learning Sparse Detector

Yen Nhi Truong Vu⋆, Dan Guo⋆, Ahmed Taha, Jason Su, Thomas Paul Matthews

WhiteRabbit.AI

**Abstract.** Deep-learning-based object detection methods show promise for improving screening mammography, but high rates of false positives can hinder their effectiveness in clinical practice. To reduce false positives, we identify three challenges: (1) unlike natural images, a malignant mammogram typically contains only one malignant finding; (2) mammography exams contain two views of each breast, and both views ought to be considered to make a correct assessment; (3) most mammograms are negative and do not contain any findings. In this work, we tackle the three aforementioned challenges by: (1) leveraging Sparse R-CNN and showing that sparse detectors are more appropriate than dense detectors for mammography; (2) including a multi-view cross-attention module to synthesize information from different views; (3) incorporating multi-instance learning (MIL) to train with unannotated images and perform breast-level classification. The resulting model, M&M, is a **M**ulti-view and **M**ulti-instance learning system that can both localize malignant findings and provide breast-level predictions. We validate M&M's detection and classification performance using five mammography datasets. In addition, we demonstrate the effectiveness of each proposed component through comprehensive ablation studies.
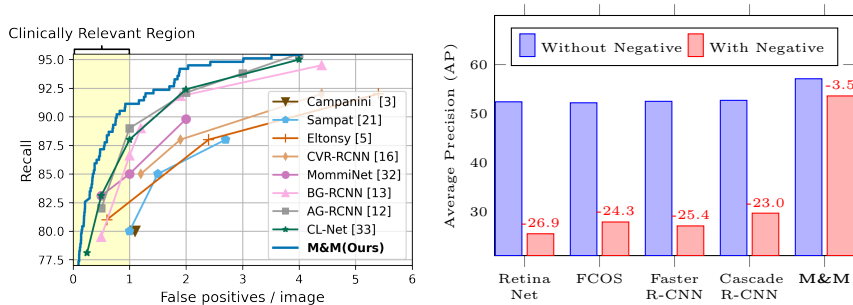
**Keywords:** Mammography · Detection · Classification · False positive

## 1 Introduction

Screening mammography helps detect breast cancer earlier and has reduced the breast cancer mortality rate significantly [4]. Computer-aided diagnosis (CAD) software was developed to aid radiologists, but its effectiveness has been questioned following recent large-scale clinical studies [6]. In particular, the high rate of false positive (FP) predictions of CAD can cause a significant reduction in radiologists' specificity [6]. Surprisingly, recent deep learning literature [3, 5, 13, 16, 20, 21, 32] focuses on improving recall without considering the need to operate at low FP rates. As shown in Fig. 1a, most works focus on reporting recalls outside the clinically relevant region of less than 1 FP/image.

---

⋆ Equal Contribution

(a) Free response operating characteristic (FROC) curves on DDSM.

(b) Quantitative detection evaluation with and without negative images on OPTIMAM.

Fig. 1: Two gaps between deep learning literature and clinical applicability. **(a)** Few works report detailed performance in the clinically relevant region of less than 1 FP/image. M&M surpasses previous works by a large margin in this region. **(b)** Typical evaluation datasets are not representative: they contain from zero (CBIS-DDSM [9]) to few negative cases (DDSM [8], INBreast [17]). To illustrate the distribution shift, we train four popular dense detectors using a standard setup that includes only annotated malignant and benign cases [1,13, 16]. We utilize OPTIMAM [7], a large dataset with a significant proportion of negatives (Tab. 1), for training and evaluation. Across all dense models, there is a large performance drop in the clinically representative setting that includes negative images. This means that the dense models are producing too many FPs on negative images. Our model, M&M, successfully tackles this performance gap.

To tackle the high rate of false positives in mammography, we identify three challenges: (1) A malignant mammogram typically contains only one malignant finding. This is different from natural images: for example, an image in COCO contains on average 7.7 objects [11]. This calls into question the usage of dense detectors for mammography; (2) A standard screening exam consists of two views per breast. Both views are essential in making a clinical decision because a finding may appear suspicious in one view but not the other; (3) Most mammograms are negative: they do not contain any findings. However, excluding negative images from training and evaluation leads to a distribution shift since negative images are abundant in clinical practice. Concretely, the false positive rate is low for a typical evaluation data distribution but much higher for a clinically-representative data distribution, as shown in Fig. 1b.

In this work, we tackle these challenges and propose a **M**ulti-view and **M**ulti-instance learning system, **M&M**. M&M is an end-to-end system that detects malignant findings and provides breast-level classification. To achieve these goals, M&M leverages three components: (1) Sparse R-CNN to replace dense anchors with a set of sparse proposals; (2) Multi-view cross-attention to synthesize information from two views and iteratively refine the predictions, and (3) Multi-

instance learning (MIL) to include negative images during training. Ultimately, each component contributes to our goal of reducing false positives.

We validate M&M through evaluation on five datasets: two in-house datasets, two public datasets — DDSM [8] and CBIS-DDSM [9], and OPTIMAM [7]. We perform ablation studies to verify the contribution of each component of M&M. To summarize, our contributions are:

1. We show that sparsity of proposals is beneficial to the analysis of mammograms, which have low disease prevalence (Sec. 2.1). With Sparse R-CNN, M&M generalizes better to clinically-representative data, where the majority of images are negative, *i.e.*, have no findings (Tab. 2);
2. We incorporate a simple and efficient cross-view multi-head attention module for mammography analysis (Sec. 2.2). With multi-view reasoning, M&M improves the recall at 0.1 FP/image by 8.6%, as shown in Fig. 4;
3. We leverage MIL to include images without bounding boxes during training (Sec. 2.3). Accordingly, M&M sees seven times more images during training. With MIL, M&M improves the recall at 0.1 FP/image by 12.6% (Fig. 4). Furthermore, M&M can provide breast-level classification predictions, achieving AUCs of more than 0.88 on four different datasets (Tab. 3).

## 2   M&M: A Multi-view and MIL System

### 2.1   Sparse R-CNN with Dual Classification Heads

The sparsity of malignant findings calls into question the use of dense detectors. As shown in Fig. 1b, dense detectors generalize poorly to negative images as they produce too many false positives. Thus, we propose to use Sparse R-CNN [24].

Sparse R-CNN utilizes a sparse set of $N$ learnable proposals consisting of $\mathbf{b_0} \in \mathbb{R}^{N \times 4}$ coordinates and $\mathbf{h_0} \in \mathbb{R}^{N \times D}$ features. The architecture uses 6 cascading heads to iteratively refine the proposals. Within the $i^{\text{th}}$ head, the proposals $\mathbf{h}_{i-1}$ first interact with themselves via self-attention, and then generate DynamicConv (Fig. 4, [24]) to interact with RoI features cropped by $\mathbf{b}_{i-1}$. The resulting outputs $\mathbf{h}_i \in \mathbb{R}^{N \times D}$ are features for the $(i+1)^{\text{th}}$ head. In addition, a regression module is applied to $\mathbf{h}_i$ to generate boxes $\mathbf{b}_i \in \mathbb{R}^{N \times 4}$, and a classification module generates scores $\mathbf{p}_i \in \mathbb{R}^{N \times C}$, with $C$ being the number of classes.

We modify Sparse R-CNN to include dual classification modules (Fig. 2). First, an objectness module produces objectness logits $\mathbf{o}_i \in \mathbb{R}^N$ to distinguish all findings — malignant and benign — from the background. By utilizing all findings, the objectness head increases the training sample size [1,13,16], but also increases FPs because it flags benign findings. To mitigate this side effect, we include a dedicated malignancy module $[\mathbf{W}_i, \mathbf{b}_i]$ to generate malignancy logits $\mathbf{m}_i \in \mathbb{R}^N$ that is trained to distinguish malignant from benign findings:

$$\mathbf{m}_i = \mathbf{o}_i - \text{SoftPlus}(\mathbf{W}_i \mathbf{h}_i + \mathbf{b}_i). \tag{1}$$

The strictly positive function $\text{SoftPlus}(x) = \log(1 + e^x)$ is chosen to enforce consistency: a high objectness logit $\mathbf{o}_i$ is required to generate a high malignancy
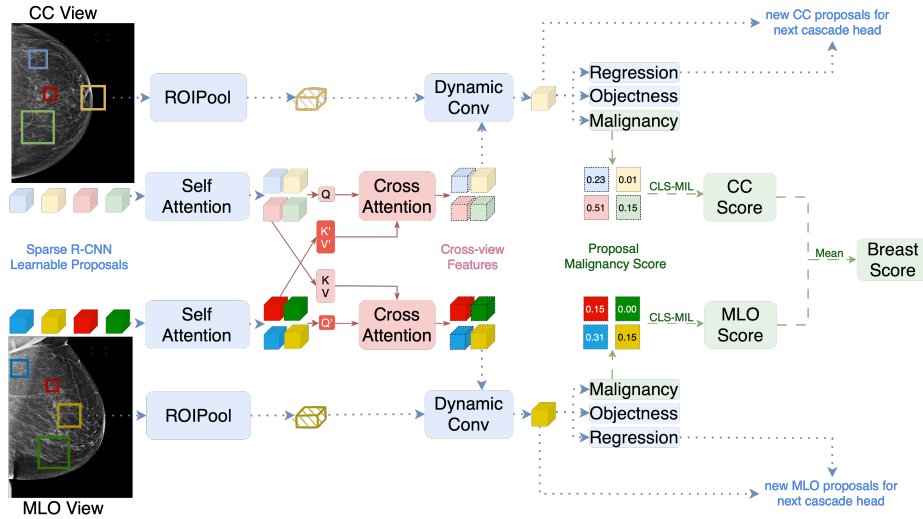
Fig. 2: M&M tackles false positives through (1, blue, dotted arrows) leveraging the Sparse R-CNN cascade architecture to iteratively refine sparse learnable proposals into predictions, (2, red, solid arrows) incorporating a cross-attention module to reason about relations between objects across two views, and (3, green, dashed arrows) utilizing image and breast MIL pooling to train with images that do not have lesion annotations.

logit $\mathbf{m}_i$. Thus, at the finding level, we obtain the following loss

$$\mathcal{L}_{\text{lesion}} = \mathcal{L}_{\text{malignant}} + \mathcal{L}_{\text{objectness}} + 2\mathcal{L}_{\text{giou}} + 5\mathcal{L}_{\text{L1}}, \tag{2}$$

where $\mathcal{L}_{\text{giou}}$ and $\mathcal{L}_{\text{L1}}$ are regression losses as in Sparse R-CNN. $\mathcal{L}_{\text{objectness}}$ and $\mathcal{L}_{\text{malignancy}}$ are focal losses applied to the predicted objectness $\mathbf{o}_i$ and the predicted malignancy $\mathbf{m}_i$ across all cascading heads $1 \leq i \leq 6$, respectively.

## 2.2   Multi-view Reasoning

A standard screening exam includes two standard views of each breast. The craniocaudal (CC) view is taken from the top down, while the mediolateral oblique (MLO) view is captured from the side at an oblique angle. Radiologists examine both views when making a clinical decision as a finding may look innocuous in one view but suspicious in the other.

To enable multi-view reasoning, M&M incorporates a cross-attention module [28] into every cascading head. Recall that within the $i^{\text{th}}$ cascading head, self-attention is first applied to proposal features $\mathbf{h}_{i-1}$ to reason about the relations between objects. After this self-attention module, we introduce a cross-attention module (Fig. 2, Appendix Algo. 1) to reason about the relations between CC

view feature $\mathbf{h}_{i-1}^{\mathrm{CC}}$ and MLO view feature $\mathbf{h}_{i-1}^{\mathrm{MLO}}$:

$$\tilde{\mathbf{h}}_{i-1}^{\mathrm{CC}} = \mathbf{h}_{i-1}^{\mathrm{CC}} + \mathrm{MultiHeadAttn}(Q = \mathbf{h}_{i-1}^{\mathrm{CC}}, V = \mathbf{h}_{i-1}^{\mathrm{MLO}}, K = \mathbf{h}_{i-1}^{\mathrm{MLO}}), \quad (3)$$

$$\tilde{\mathbf{h}}_{i-1}^{\mathrm{MLO}} = \mathbf{h}_{i-1}^{\mathrm{MLO}} + \mathrm{MultiHeadAttn}(Q = \mathbf{h}_{i-1}^{\mathrm{MLO}}, V = \mathbf{h}_{i-1}^{\mathrm{CC}}, K = \mathbf{h}_{i-1}^{\mathrm{CC}}). \quad (4)$$

The enhanced embeddings $\tilde{\mathbf{h}}_{i-1}^{\mathrm{CC}}$, $\tilde{\mathbf{h}}_{i-1}^{\mathrm{MLO}}$ then generate DynamicConv to interact with RoI features and produce new features $\mathbf{h}_i^{\mathrm{CC}}$, $\mathbf{h}_i^{\mathrm{MLO}}$ for the $(i+1)^{\mathrm{th}}$ head. Thus, with the proposed cross-attention module, the CC view's proposal features are refined iteratively using the MLO view's proposal features and vice versa.

### 2.3   Multi-instance Learning

Mammogram annotation is costly to obtain due to a dependency on radiologists. This high cost means that bounding boxes are often unavailable. Further, most mammograms are negative: they do not contain any findings. Yet, a model generalizes poorly if these negative images are dropped during training (Fig. 1b).

   Since image- and breast-level labels are available, we adopt an MIL module to include images without bounding boxes during training. To compute image- and breast-level scores, we leverage the proposal malignancy logits $\mathbf{m}_i$ (Eq. (1)). Since an image is malignant if it contains a malignant lesion, we obtain image-level scores by applying the NoisyOR function $f(\mathbf{x}) = 1 - \prod_{k=1}^{N}(1 - \mathbf{x}[k])$ to the malignancy probabilities $\mathbf{p}_i = \mathrm{Sigmoid}(\mathbf{m}_i) \in \mathbb{R}^N$. Next, as CC and MLO views offer complimentary information on a breast, we obtain breast-level malignancy score by averaging the image-level scores across these views.

   We apply cross-entropy losses $\mathcal{L}_{\mathrm{image}}$ and $\mathcal{L}_{\mathrm{breast}}$ at the image and breast level for all training samples. The lesion loss $\mathcal{L}_{\mathrm{lesion}}$ (Eq. (2)) is only applied for annotated lesions. We thus obtain the following total training loss for M&M:

$$\mathcal{L} = \mathbb{1}_{\mathrm{annotated\ lesion}}\mathcal{L}_{\mathrm{lesion}} + 0.5\mathcal{L}_{\mathrm{image}} + 0.5\mathcal{L}_{\mathrm{breast}}. \quad (5)$$

## 3   Experiments

**Implementation Details.** We use PyTorch 1.10. The training settings follow Sparse R-CNN [24]. We apply random horizontal flipping and random rotation. We resize the images' shorter edges to 2560 with the larger edges no longer than 3328. We utilize a COCO-pretrained PVT-B2-Li backbone [30]. We use AdamW optimizer with $5 \times 10^{-5}$ learning rate and 0.0001 weight decay. The model is trained for 9000 iterations, and the learning rate is scaled by 0.1 at the 6750 and 8250 iterations. Each batch contains 16 breasts (32 images). We employ a 1:1 sampling ratio between unannotated and annotated images.

**Datasets.** We utilize three 2D digital mammography datasets: (1) *OPTIMAM*: a development dataset derived from the OPTIMAM database [7], which is funded by Cancer Research UK. We split the data into train/val/test with an 80:10:10

Table 1: Dataset statistics. We report the number of breasts in each dataset, broken down by 3 categories: malignant, benign, and negative. Malignant breasts contain findings with positive biopsy outcomes. Benign breasts contain findings that are determined to be non-malignant after additional follow-up. Negative breasts do not contain any radiologist-marked findings. In the parentheses, we report the number of breasts with bounding box annotations. "Bbox" indicates whether bounding box annotations are available.

| Datasets | Bbox | Malignant (Ann.) | | Benign (Ann.) | | Negative (Ann.) | |
|---|---|---|---|---|---|---|---|
| OPTIMAM | ✓ | 4,838 | (4,245) | 1,999 | (567) | 26,003 | (2) |
| Inhouse-A | | 496 | (0) | 2,128 | (0) | 2074 | (0) |
| Inhouse-B | | 243 | (0) | 7,797 | (0) | 47,929 | (0) |
| DDSM | ✓ | 624 | (624) | 555 | (555) | 2,877 | (1) |
| CBIS-DDSM | ✓ | 312 | (310) | 347 | (336) | 0 | (0) |

Table 2: Quantitative detection evaluation on OPTIMAM. $\Delta$ denotes the AP gap between evaluating with and without negative images.

| Model | $AP_{mb}$ | AP | $\Delta$ | R@0.1 | R@0.25 | R@0.5 |
|---|---|---|---|---|---|---|
| RetinaNet [10] | 52.4 | 25.5 | -26.9 | 53.3 | 73.1 | 83.0 |
| FCOS [26] | 52.2 | 27.9 | -24.3 | 52.0 | 77.4 | 87.0 |
| Faster R-CNN [19] | 52.5 | 27.1 | -25.4 | 51.5 | 71.2 | 84.1 |
| Cascade R-CNN [2] | 52.7 | 29.7 | -23.0 | 54.9 | 77.0 | 86.2 |
| Sparse R-CNN [24] | 53.2 | 36.2 | -17.0 | 64.3 | 77.0 | 85.5 |
| **M&M (ours)** | **57.1** | **53.6** | **-3.5** | **87.7** | **90.9** | **92.5** |

ratio at the patient level; (2) *Inhouse-A*: an evaluation dataset collected from a U.S. multi-site mammography operator; (3) *Inhouse-B*: an evaluation dataset collected from a U.S. academic hospital (see [18], Sec. 2.2 for more details on the inhouse datasets). We also utilize two film mammography datasets: (4) *DDSM*: a dataset maintained at the University of South Florida [8]. We followed the methods by [3,5,13,16] to split the test set; (5) *CBIS-DDSM:* a curated subset of DDSM [9]. We only include breasts that have one CC view and one MLO view. Dataset statistics are reported in Tab. 1.

**Metrics.** We report average precision with Intersection over Union from 0.25 to 0.75. $AP_{mb}$ denotes average precision on the set of annotated malignant and benign images. AP denotes average precision when all data is included. We report free response operating characteristic (FROC) curves and recalls at various FP/image (R@t). Following [3,5,16,29], a proposal is considered true positive if its center lies within the ground truth box. For classification, we report the area under the receiver operating characteristic curve (AUC).

**Detection Results.** Tab. 2 presents quantitative detection evaluation on OP-TIMAM. All dense detectors [2,10,19,26] suffer a large $\Delta$ gap of more than 23 points (pt) between excluding and including negative images. Large $\Delta$ means the models are producing too many FPs on negative images. Sparse R-CNN [24]
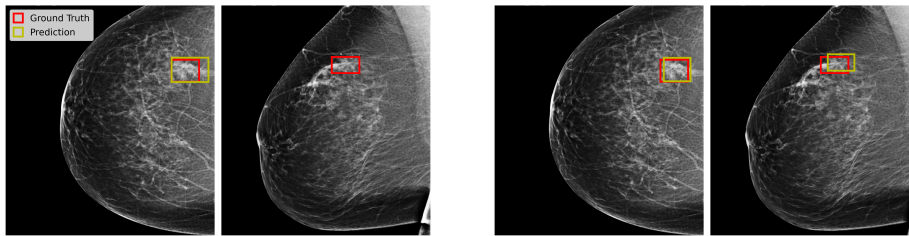
Fig. 3: Qualitative Evaluation. **Left**: Model without multi-view (row 4 of Fig. 4) produces a loose box on the CC view and misses the finding on the MLO view. **Right**: M&M produces tight boxes around the finding in both views.

Table 3: Quantitative classification evaluation. (a) On three private datasets, we use two open-sourced mammography classifiers as baselines [23, 25]. All models were trained only on OPTIMAM. We report AUC at both the breast and the exam level, except for Inhouse-A, where breast-level labels are unavailable. (b) We train M&M on CBIS-DDSM and compare breast AUC with recent literature. (* Tulder *et al.* [27] report results using five-fold cross validation.)

(a) Private Datasets

| Model | OPTIMAM | | Inhouse-A | Inhouse-B | |
|---|---|---|---|---|---|
| | Breast AUC | Exam AUC | Exam AUC | Breast AUC | Exam AUC |
| GMIC [23] | 0.911 | 0.896 | 0.814 | 0.815 | 0.796 |
| HCT [25] | 0.923 | 0.912 | 0.816 | 0.817 | 0.793 |
| **M&M (ours)** | **0.960** | **0.942** | **0.920** | **0.910** | **0.898** |

(b) CBIS-DDSM

| Model | Breast AUC |
|---|---|
| ResNet50 [14] | 0.724 |
| Shared ResNet [31] | 0.735 |
| PHResNet50 [14] | 0.739 |
| Cross-view Transformer [27] | 0.803* |
| **M&M (ours)** | **0.883** |

generalizes significantly better with a gap of 17pt. This shows the importance of sparsity for reducing FP. By adding both multi-view and MIL, M&M successfully reduces the $\Delta$ gap to 3.5pt. With this performance gap closed, M&M is able to achieve a high recall of 87.7% at just 0.1 FP/image.

Fig. 1a compares M&M with recent literature evaluated on DDSM. M&M adopts the same DDSM splits used by [3, 12, 13, 16, 33], while [5, 21, 32] use other splits. M&M (87% R@0.5) outperforms all recent SOTA with the same test split, including 2022 SOTA [33] (83% R@0.5), by at least 4%.

**Classification Results.** Tab. 3a reports M&M's breast-level and exam-level classification results on OPTIMAM and the two inhouse datasets. We use GMIC [23] and HCT [25] as baselines since they are open-sourced classifiers developed for mammography. All three models were trained only on OPTIMAM. For all models, the breast-level score is the average of the CC score and MLO score, while the exam-level score is the max of the left breast score and right breast score. Both baseline models suffer large generalization drops of approximately
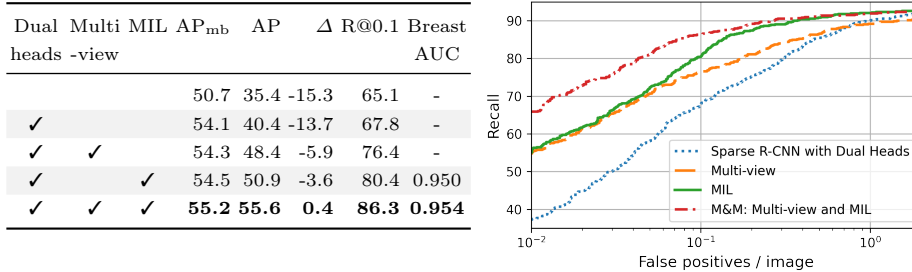
| Dual heads | Multi -view | MIL | $AP_{mb}$ | AP | $\Delta$ | R@0.1 | Breast AUC |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | 50.7 | 35.4 | -15.3 | 65.1 | - |
| ✓ | | | 54.1 | 40.4 | -13.7 | 67.8 | - |
| ✓ | ✓ | | 54.3 | 48.4 | -5.9 | 76.4 | - |
| ✓ | | ✓ | 54.5 | 50.9 | -3.6 | 80.4 | 0.950 |
| **✓** | **✓** | **✓** | **55.2** | **55.6** | **0.4** | **86.3** | **0.954** |

Fig. 4: Effect of M&M's components on classification and detection performance.

0.08–0.12 exam AUC when evaluated on Inhouse-A and Inhouse-B. In comparison, M&M has smaller performance gaps of 0.02 on Inhouse-A and 0.04 on Inhouse-B. Similar observations for other classifiers, such as EfficientNet, are reported in the appendix.

Tab. 3b compares M&M with recent literature reporting on the public CBIS-DDSM dataset. In particular, M&M outperforms the cross-view transformer [27] and PHResNet50 [14] by 0.08 and 0.14 breast AUC, respectively.

**Qualitative Evaluation.** Fig. 3 presents a qualitative evaluation of the multi-view module. With multi-view, M&M produces a tighter box on the CC view and recovers a missed finding on the MLO view.

**Ablation Studies.** Fig. 4 presents ablation results using the OPTIMAM validation split. On the left, we demonstrate how each component of M&M contributes to closing the gap $\Delta$ between evaluating with and without negative images. Notably, without using any extra training samples, multi-view reasoning reduces $\Delta$ to only −5.9pt (Row 3). MIL allows the model to train with significantly more negative images, reducing $\Delta$ to −3.6pt (Row 4). On the right of Fig. 4, the FROC curves show how each component of M&M improves recall significantly at low FP/image. In particular, M&M's recall at 0.1FP/image is 86.3%, +21.2% over vanilla Sparse R-CNN.

**Further studies.** In the appendix, we present more qualitative evaluation as well as further ablation studies on (1) number of learnable proposals, (2) different MIL schemes, (3) backbone choices and (4) positional encoding.

## 4    Discussion and Conclusion

We present M&M, an end-to-end model leveraging multi-view reasoning and multi-instance learning for mammography detection and classification.

As a detector, M&M offers significant improvement in recall at low FP/image (Fig. 1a, Tab. 2). This success comes from three points of advancement. First, unlike previous works that do not consider the impact of sparsity [13,16,33], we show that sparsity of proposals is beneficial for false positive reduction (Tab. 2). Second, M&M incorporates multi-view reasoning through iterative application of cross-attention and proposal refinement in the cascading heads. M&M's multi-view module is effective (Fig. 4) yet simple, requiring neither positional encoding [13,16,32] nor extra proposal correspondence annotations [33]. Finally, our MIL formulation allows for training with representative data distribution in an end-to-end one stage pipeline. This is more advantageous than previous pipelines that require additional stages or classifiers to reduce false positives [15,22,29].

As a classifier, M&M establishes strong performance on several datasets (Tab. 3). M&M offers two advantages over image classifiers: (1) Image classifiers are often pre-trained as patch classifiers with patches cropped from bounding box annotations [14,23,25]. In comparison, M&M utilizes these bounding boxes to learn localization and can be trained directly in a single stage from COCO/ImageNet weights; (2) Image classifiers offer limited explainability, while M&M's breast-level prediction is more interpretable through its localization ability.

# References

1. Agarwal, R., Diaz, O., Lladó, X., Yap, M.H., Martí, R.: Automatic mass detection in mammograms using deep convolutional neural networks. Journal of Medical Imaging **6**(3), 031409 (2019) 2, 3
2. Cai, Z., Vasconcelos, N.: Cascade r-cnn: Delving into high quality object detection. In: CVPR (2018) 6
3. Campanini, R., Dongiovanni, D., Iampieri, E., Lanconelli, N., Masotti, M., Palermo, G., Riccardi, A., Roffilli, M.: A novel featureless approach to mass detection in digital mammograms based on support vector machines. Physics in Medicine & Biology **49**(6), 961 (2004) 1, 6, 7
4. Duffy, S.W., Tabár, L., Yen, A.M.F., Dean, P.B., Smith, R.A., Jonsson, H., Törnberg, S., Chen, S.L.S., Chiu, S.Y.H., Fann, J.C.Y., et al.: Mammography screening reduces rates of advanced and fatal breast cancers: Results in 549,091 women. Cancer **126**(13), 2971–2979 (2020) 1
5. Eltonsy, N.H., Tourassi, G.D., Elmaghraby, A.S.: A concentric morphology model for the detection of masses in mammography. IEEE transactions on medical imaging **26**(6), 880–889 (2007) 1, 6, 7
6. Fenton, J.J., Abraham, L., Taplin, S.H., Geller, B.M., Carney, P.A., D'Orsi, C., Elmore, J.G., Barlow, W.E., Consortium, B.C.S.: Effectiveness of computer-aided detection in community mammography practice. Journal of the National Cancer institute **103**(15), 1152–1161 (2011) 1
7. Halling-Brown, M.D., Warren, L.M., Ward, D., Lewis, E., Mackenzie, A., Wallis, M.G., Wilkinson, L.S., Given-Wilson, R.M., McAvinchey, R., Young, K.C.: Optimam mammography image database: a large-scale resource of mammography images and clinical data. Radiology: Artificial Intelligence (2020) 2, 3, 5
8. Heath, M., Bowyer, K., Kopans, D., Moore, R., Kegelmeyer, P.: The digital database for screening mammography. In: Proceedings of the Fifth International Workshop on Digital Mammography. Medical Physics Publishing (2001) 2, 3, 6

9. Lee, R.S., Gimenez, F., Hoogi, A., Miyake, K.K., Gorovoy, M., Rubin, D.L.: A curated mammography data set for use in computer-aided detection and diagnosis research. Scientific data **4**(1), 1–9 (2017) 2, 3, 6

10. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV. pp. 2980–2988 (2017) 6

11. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014) 2

12. Liu, Y., Zhang, F., Chen, C., Wang, S., Wang, Y., Yu, Y.: Act like a radiologist: towards reliable multi-view correspondence reasoning for mammogram mass detection. PAMI **44**(10), 5947–5961 (2021) 7

13. Liu, Y., Zhang, F., Zhang, Q., Wang, S., Wang, Y., Yu, Y.: Cross-view correspondence reasoning based on bipartite graph convolutional network for mammogram mass detection. In: CVPR. pp. 3812–3822 (2020) 1, 2, 3, 6, 7, 9

14. Lopez, E., Grassucci, E., Valleriani, M., Comminiello, D.: Multi-view breast cancer classification via hypercomplex neural networks. arXiv:2204.05798 (2022) 7, 8, 9

15. Lotter, W., Diab, A.R., Haslam, B., Kim, J.G., Grisot, G., Wu, E., Wu, K., Onieva, J.O., Boyer, Y., Boxerman, J.L., et al.: Robust breast cancer detection in mammography and digital breast tomosynthesis using an annotation-efficient deep learning approach. Nature Medicine **27**(2), 244–249 (2021) 9

16. Ma, J., Li, X., Li, H., Wang, R., Menze, B., Zheng, W.S.: Cross-view relation networks for mammogram mass detection. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 8632–8638. IEEE (2021) 1, 2, 3, 6, 7, 9

17. Moreira, I.C., Amaral, I., Domingues, I., Cardoso, A., Cardoso, M.J., Cardoso, J.S.: Inbreast: toward a full-field digital mammographic database. Academic radiology **19**(2), 236–248 (2012) 2

18. Pedemonte, S., Tsue, T., Mombourquette, B., Vu, Y.N.T., Matthews, T., Hoil, R.M., Shah, M., Ghare, N., Zingman-Daniels, N., Holley, S., et al.: A deep learning algorithm for reducing false positives in screening mammography. arXiv preprint arXiv:2204.06671 (2022) 6

19. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. NeurIPS **28** (2015) 6

20. Ren, Y., Lu, J., Liang, Z., Grimm, L.J., Kim, C., Taylor-Cho, M., Yoon, S., Marks, J.R., Lo, J.Y.: Retina-match: Ipsilateral mammography lesion matching in a single shot detection pipeline. In: MICCAI. pp. 345–354. Springer (2021) 1

21. Sampat, M.P., Bovik, A.C., Whitman, G.J., Markey, M.K.: A model-based framework for the detection of spiculated masses on mammography a. Medical physics **35**(5), 2110–2123 (2008) 1, 7

22. Sarath, C.K., Chakravarty, A., Ghosh, N., Sarkar, T., Sethuraman, R., Sheet, D.: A two-stage multiple instance learning framework for the detection of breast cancer in mammograms. In: EMBC. IEEE (2020) 9

23. Shen, Y., Wu, N., Phang, J., Park, J., Liu, K., Tyagi, S., Heacock, L., Kim, S.G., Moy, L., Cho, K., et al.: An interpretable classifier for high-resolution breast cancer screening images utilizing weakly supervised localization. Medical image analysis **68**, 101908 (2021) 7, 9

24. Sun, P., Zhang, R., Jiang, Y., Kong, T., Xu, C., Zhan, W., Tomizuka, M., Li, L., Yuan, Z., Wang, C., et al.: Sparse r-cnn: End-to-end object detection with learnable proposals. In: CVPR (2021) 3, 5, 6, 1

25. Taha, A., Truong Vu, Y.N., Mombourquette, B., Matthews, T.P., Su, J., Singh, S.: Deep is a luxury we don't have. In: MICCAI. pp. 25–35. Springer (2022) 7, 9

26. Tian, Z., Shen, C., Chen, H., He, T.: Fcos: Fully convolutional one-stage object detection. In: CVPR (2019) 6

27. Tulder, G.v., Tong, Y., Marchiori, E.: Multi-view analysis of unregistered medical images using cross-view transformers. In: MICCAI. Springer (2021) 7, 8

28. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. NeurIPS **30** (2017) 4

29. Vu, Y.N.T., Mombourquette, B., Matthews, T.P., Su, J., Singh, S.: Wrdet: a breast cancer detector for full-field digital mammograms. In: Medical Imaging 2022: Computer-Aided Diagnosis. vol. 12033, pp. 219–230. SPIE (2022) 6, 9

30. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pvt v2: Improved baselines with pyramid vision transformer. Computational Visual Media **8**(3), 415–424 (2022) 5

31. Wu, N., Jastrzębski, S., Park, J., Moy, L., Cho, K., Geras, K.J.: Improving the ability of deep neural networks to use information from multiple views in breast cancer screening. In: Medical Imaging with Deep Learning. PMLR (2020) 7

32. Yang, Z., Cao, Z., Zhang, Y., Tang, Y., Lin, X., Ouyang, R., Wu, M., Han, M., Xiao, J., Huang, L., et al.: Momminet-v2: Mammographic multi-view mass identification networks. Medical Image Analysis **73**, 102204 (2021) 1, 7, 9

33. Zhao, Z., Wang, D., Chen, Y., Wang, Z., Wang, L.: Check and link: Pairwise lesion correspondence guides mammogram mass detection. ECCV (2022) 7, 9, 1

---

**Algorithm 1** M&M Multi-view Cross Attention. This module is to be called on L191 in the official implementation of Sparse R-CNN head.

---

```python
1  class MultiviewCrossAttn(nn.Module):
2    def __init__(self, hdim=128, nhead=8, dropout=0.1):
3      self.mv_atn = nn.MultiheadAttention(hdim, nhead)
4      self.dropout_mv = nn.Dropout(dropout)
5      self.norm_mv = nn.LayerNorm(hdim)
6    def forward(self, pro_features):
7      # collect CC and MLO features from the batch dimension
8      cc_feats = pro_features[:, :pro_features.shape[1]//2]
9      mlo_feats = pro_features[:, pro_features.shape[1]//2:]
10     # cross attention to enhance CC features
11     cc_feats2,_ = self.mv_atn(query=cc_feats, key=mlo_feats, value=mlo_feats)
12     cc_feats += self.dropout_mv(cc_feats2)
13     cc_feats = self.norm_mv(cc_feats)
14     [...] # vice versa for MLO, omitted here due to space
15     pro_features = torch.stack(cc_feats, mlo_feats, dim=1) # restacking
16     return pro_features # new features used for inst_interact (DynamicConv)
```

---

Table A1: Ablation study on (a) effect of number of learnable proposals in a multi-view only model, and (b) effect of positional encoding. Results are on OPTIMAM validation set.

(a) The gap between evaluating with and without negatives $\Delta$ worsens as $N$ increases, showing that multi-view reasoning benefits from sparsity.

| No. Proposals | Training Time | $AP_{mb}$ | AP | $\Delta$ |
|---|---|---|---|---|
| 10 | 7.5h | 53.6 | 48.4 | **-5.2** |
| 40 | 7.7h | **54.3** | **48.4** | -5.9 |
| 100 | 8.0h | 53.2 | 47.1 | -6.1 |
| 400 | 8.9h | 53.0 | 45.8 | -7.2 |

(b) Different from [13, 16, 33], we found that positional encodings deliver insignificant boosts in AP and breast AUC. Our observation is similar to Sparse R-CNN's observations ( [24], Tab. 10).

| Positional Encoding | $AP_{mb}$ | AP | Breast AUC |
|---|---|---|---|
| None | **55.2** | 55.6 | 0.954 |
| Proposal center | 55.1 | 55.2 | **0.955** |
| Nipple distance | 54.4 | **55.7** | **0.955** |

Table A2: Effect of MIL approach. We also experiment with learnable image-MIL by applying a FC layer on (1) GAP: the Global Average Pooled proposal features, and (2) CLS-token: a BERT-like token that summarizes proposal features. Results are on OPTIMAM validation set.

| Image MIL | Breast MIL | $AP_{mb}$ | AP | $\Delta$ | Breast AUC |
|---|---|---|---|---|---|
| Max | Max | 54.7 | 54.0 | -0.7 | 0.952 |
|  | Mean | 54.5 | 54.6 | 0.1 | 0.953 |
| Noisy-OR | Max | 55.1 | 55.5 | **0.4** | **0.954** |
|  | Mean | 55.2 | **55.6** | **0.4** | **0.954** |
| GAP | Max | 55.9 | 54.6 | -1.3 | 0.954 |
|  | Mean | 56.1 | 54.0 | -2.1 | 0.949 |
| CLS-token | Max | 55.3 | 53.7 | -1.6 | 0.947 |
|  | Mean | 56.2 | 55.0 | -1.2 | 0.951 |

Table A3: Quantitative detection evaluation with different backbones on two test sets. On OPTIMAM, across all different backbones, M&M has a small gap $\Delta$ between evaluating with and without negative images. On DDSM, M&M achieves more than 83% recall at 0.5 FP/image across three different backbones.

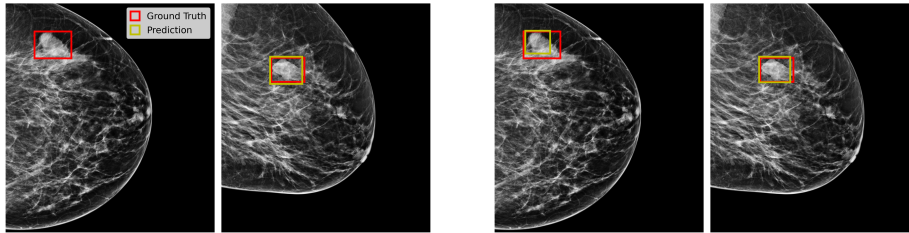| Dataset | Backbone | $AP_{mb}$ | AP | $\Delta$ | R@0.1 | R@0.25 | R@0.5 |
|---------|----------|-----------|------|----------|-------|--------|-------|
| OPTIMAM | GMIC | 50.7 | 44.6 | -6.1 | 75.4 | 82.5 | 86.5 |
| | EfficientNet-B0 | 51.9 | 45.4 | -6.5 | 78.8 | 87.1 | 89.7 |
| | ResNet-50 | 52.8 | 47.0 | -5.8 | 80.4 | 87.2 | 90.0 |
| | PVT-B2-Li | **57.1** | **53.6** | **-3.5** | **87.7** | **90.9** | **92.5** |
| DDSM | GMIC | 31.1 | 27.0 | -4.1 | 52.9 | 72.2 | 79.2 |
| | EfficientNet-B0 | 39.1 | 35.2 | -3.9 | 66.4 | 74.6 | 83.2 |
| | ResNet-50 | 38.9 | 36.6 | -2.3 | 75.1 | 79.2 | 83.5 |
| | PVT-B2-Li | **39.2** | **37.0** | **-2.2** | **80.4** | **82.6** | **87.2** |



Fig. A1: Additional Qualitative Evaluation. **Left**: without multi-view, the model misses a mass on the CC view even though it was able to detect the mass on the MLO view. **Right**: with multi-view, M&M recalls the mass on both views.

Table A4: Quantitative classification evaluation with different backbones on 3 datasets. All models are trained using OPTIMAM. M&M column denotes whether the model was a classifier (-) or M&M with the row's backbone (✓).

| Backbone | M&M | OPTIMAM | | Inhouse-A | Inhouse-B | |
|----------|-----|---------|---|-----------|-----------|---|
| | | Breast AUC | Exam AUC | Exam AUC | Breast AUC | Exam AUC |
| GMIC | - | 0.911 | 0.896 | 0.814 | 0.815 | 0.796 |
| | ✓ | 0.920 | 0.900 | 0.835 | 0.843 | 0.822 |
| EfficientNet-B0 | - | 0.940 | 0.922 | 0.787 | 0.850 | 0.826 |
| | ✓ | 0.941 | 0.913 | 0.840 | 0.852 | 0.832 |
| PVT-B2-Li | - | 0.949 | 0.933 | 0.820 | 0.867 | 0.846 |
| | ✓ | **0.960** | **0.942** | **0.920** | **0.910** | **0.898** |

Table A5: Quantitative classification evaluation with different backbones on CBIS-DDSM. All models are trained using CBIS-DDSM.

| Metric/ Backbone | GMIC | EfficientNet-B0 | ResNet-50 | PVT-B2-Li |
|------------------|------|-----------------|-----------|-----------|
| Breast AUC | 0.839 | 0.865 | 0.836 | 0.883 |
| Exam AUC | 0.835 | 0.865 | 0.829 | 0.883 |