# Anatomy-Aware Siamese Network: Exploiting Semantic Asymmetry for Accurate Pelvic Fracture Detection in X-ray Images

Haomin Chen[⋆1,2], Yirui Wang[⋆1], Kang Zheng[1], Weijian Li[3], Chi-Tung Chang[5], Adam P. Harrison[1], Jing Xiao[4], Gregory D. Hager[2], Le Lu[1], Chien-Hung Liao[5], and Shun Miao[1]

[1] PAII Inc., Bethesda, MD, USA
[2] Departemnt of Computer Science, Johns Hopkins University, Baltimore, MD, USA
[3] Department of Computer Science, University of Rochester, NY, USA
[4] Ping An Technology, Shenzhen, China
[5] Chang Gung Memorial Hospital, Linkou, Taiwan, ROC

**Abstract.** Visual cues of enforcing bilaterally symmetric anatomies as normal findings are widely used in clinical practice to disambiguate subtle abnormalities from medical images. So far, inadequate research attention has been received on effectively emulating this practice in CAD methods. In this work, we exploit semantic anatomical symmetry or asymmetry analysis in a complex CAD scenario, i.e., anterior pelvic fracture detection in trauma PXRs, where semantically pathological (refer to as fracture) and non-pathological (*e.g.* pose) asymmetries both occur. Visually subtle yet pathologically critical fracture sites can be missed even by experienced clinicians, when limited diagnosis time is permitted in emergency care. We propose a novel fracture detection framework that builds upon a Siamese network enhanced with a spatial transformer layer to holistically analyze symmetric image features. Image features are spatially formatted to encode bilaterally symmetric anatomies. A new contrastive feature learning component in our Siamese network is designed to optimize the deep image features being more salient corresponding to the underlying semantic asymmetries (caused by pelvic fracture occurrences). Our proposed method have been extensively evaluated on 2,359 PXRs from unique patients (the largest study to-date), and report an area under ROC curve score of 0.9771. This is the highest among state-of-the-art fracture detection methods, with improved clinical indications.

**Keywords:** Anatomy-Aware Siamese Network, Semantic Asymmetry, Fracture Detection, X-ray Images

## 1 Introduction

The CAD of abnormalities in medical images is among the most promising applications of computer vision in healthcare. In particular, X-ray CAD represents
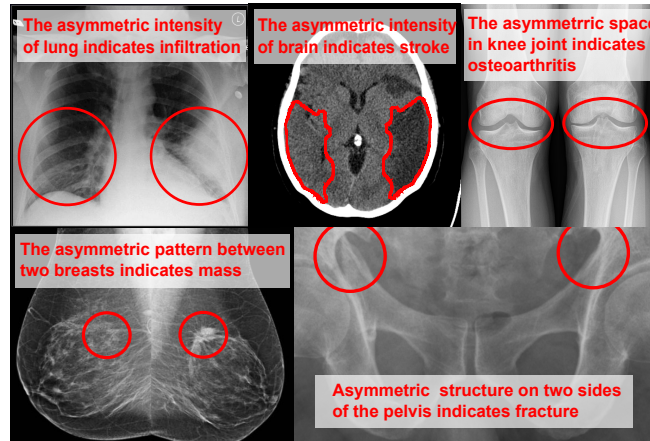
---

[⋆] equal contribution

**Fig. 1.** Example medical images where anatomical symmetry helps to detect abnormalities. The top 3 images represents infiltration in chest X-Rays, stroke in brain CT, and osteoarthritis in knee X-Rays. The bottom 2 images represent masses in mammography and fractures in PXRs. These abnormalities can be better differentiated when the anatomically symmetric body parts are compared.

an important research focus [5,34,28,20,15,4,25]. However, the high variations of abnormalities in medical imagery pose non-trivial challenges in differentiating pathological abnormalities from radiological patterns caused by normal anatomical and imaging-condition differences. At the same time, many anatomical structures are bilaterally symmetric (e.g., the brain, skeleton and breast) which suggests that the detection of abnormal radiological findings can exploit semantically symmetric anatomical regions (Figure 1). Indeed, using bilaterally symmetric visual cues to confirm suspicious findings is a strongly recommended and widely adopted clinical practice [7]. Our aim is to emulate this practice in CAD and apply it to the problem of effectively detecting subtle but critical anterior pelvic fractures in trauma PXRs.

Several studies have investigated the use of symmetry cues for CAD, aiming to find abnormalities in brain structures in neuro-imaging [32,18,22], breasts in mammograms [24], and stroke in CT [1]. All of these works directly employ symmetry defined on the image or shape space. However, under less constrained scenarios, especially the ones using projection-based imaging modalities in an emergency room setting, *e.g.*, PXRs, image asymmetries do not always indicate positive clinical findings, as they are often caused by other non-pathological factors like patient pose, bowel gas patterns, and clothing. For these settings, a workflow better mirroring the clinical practice, *i.e.* robust analysis across semantic *anatomical* symmetries, is needed. Using semantic anatomical symmetry to facilitate CAD in such complex scenarios has yet to be explored.

To bridge this gap, we propose an AASN to effectively exploit semantic anatomical symmetry in complex imaging scenarios. Our motivation comes from the detection of pelvic fractures in emergency-room PXRs. Pelvic fractures are among the most dangerous and lethal traumas, due to their high association with massive internal bleeding. Non-displaced fractures, *i.e.*, fractures that cause no displacement of the bone structures, can be extraordinarily difficult to detect, even for experienced clinicians. Therefore, the combination of difficult detection coupled with extreme and highly-consequential demands on performance motivates even more progress. Using anatomical symmetry to push the performance even higher is a critical gap to fill.

In AASNs, we employ fully convolutional Siamese networks [11] as the backbone of our method. First, we exploit symmetry cues by anatomically reparameterizing the image using a powerful graph-based landmark detection [21]. This allows us to create an anatomically-grounded warp from one side of the pelvis to the other. While previous symmetry modeling methods rely on image-based spatial alignment before encoding [24], we take a different approach and perform feature alignment after encoding using a spatial transformer layer. This is motivated by the observation that image *asymmetry* in PXRs can be caused by many factors, including imaging angle and patient pose. Thus, directly warping images is prone to introducing artifacts, which can alter pathological image patterns and make them harder to detect. Since image asymmetry can be semantically pathological, *i.e.*, fractures, and non-pathological, *e.g.*, imaging angle and patient pose, we propose a new contrastive learning component in Siamese network to optimize the deep image features being more salient corresponding to the underlying semantic asymmetries (caused by fracture). Crucially, this mitigates the impact of distracting asymmetries that may mislead the model. With a sensible embedding in place, corresponding anatomical regions are jointly decoded for fracture detection, allowing the decoder to reliably discover fracture-causing discrepancies.

In summary, our main contributions are four folds.

- We present a clinically-inspired (or reader-inspired) and computationally principled framework, named AASN, which is capable of effectively exploiting anatomical landmarks for semantic asymmetry analysis from encoded deep image features. This facilitates a high performance CAD system of detecting both visually evident and subtle pelvic fractures in PXRs.
- We systematically explore plausible means for fusing the image based anatomical symmetric information. A novel Siamese feature alignment via spatial transformer layer is proposed to address the potential image distortion drawback in the prior work [24].
- We describe and employ a new contrastive learning component to improve the deep image feature's representation and saliency reflected from semantically pathological asymmetries. This better disambiguates against the existing visual asymmetries caused by non-pathological reasons.
- Extensive evaluation on real clinical dataset of 2,359 PXRs from unique patients is conducted. Our results show that AASN simultaneously increases

the AUC and the average precision from 96.52% to 97.71% or from 94.52% to 96.50%, respectively, compared to a strong baseline model that does not exploit symmetry or asymmetry. *More significantly, the pelvic fracture detection sensitivity or recall value has been boosted from 70.79% to 77.87% when controlling the false positive (FP) rate at 1%.*

## 2    Related Work

**Computer-Aided Detection and Diagnosis in Medical Imaging.** In recent years, motivated by the availability of public X-ray datasets, X-ray CAD has received extensive research attention. Many works have studied abnormality detection in CXRs [5,34,28,20]. CAD of fractures in musculoskeletal radiographs is another well studied field [6,8,35]. Since many public X-ray datasets only have image-level labels, many methods formulate abnormality detection as an image classification problem and use class activation maps [38] for localization [28,34]. While abnormalities that involve a large image area (e.g., atelectasis, cardiomegaly) may be suitable for detection via image classification, more localized abnormalities like masses and fractures are in general more difficult to detect without localization annotations. While methods avoiding such annotations have been developed [20,35], we take a different approach and use point-based localizations for annotations, which are minimally laborious and a natural fit for ill-defined fractures. Another complementary strategy to improve abnormality detection is to use anatomical and pathological knowledge and heuristics to help draw diagnostic inferences [23]. This is also an approach we take, exploiting the bilateral symmetry priors of anatomical structures to push forward classification performance.

**Image based Symmetric Modeling for CAD.** Because many human anatomies are left-right symmetric (e.g., brain, breast, bone), anatomical symmetry has been studied for CAD. The shape asymmetry of subcortical brain structures is known to be associated with Alzheimer's disease and has been measured using both analytical shape analysis [32,18] and machine learning techniques [22]. A few attempts have been explored using symmetric body parts for CAD [1,24]. For instance, Siamese networks [11] have been used to combine features of the left and right half of brain CTs for detecting strokes. A Siamese Faster-RCNN approach was also proposed to detect masses from mammograms by jointly analyzing left and right breasts [24]. Yet, existing methods directly associate asymmetries in the image space with pathological abnormalities. While this assumption may hold in strictly controlled imaging scenarios, like brain CT/MRIs and mammograms, this rarely holds in PXRs, where additional asymmetry causing factors are legion, motivating the more anatomically-derived approach to symmetry that we take.

**Siamese Network and Contrastive Learning.** Siamese networks are an effective method for contrastive learning that uses contrastive loss to embed semantically similar samples closer together and dissimilar images further away [11]. Local similarities have also been learned using Siamese networks [37]
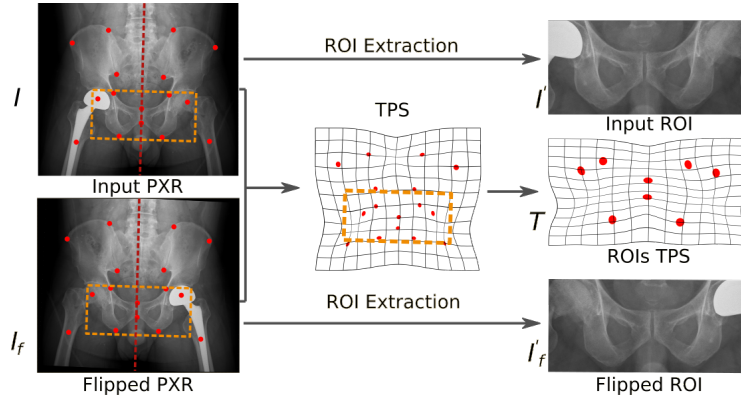
**Fig. 2.** Illustration of ROI and warp generation steps.

and applied to achieve image matching/registration [26,29]. The embedding learned by Siamese networks has also been applied to one-shot image recognition [17] and human re-identification [31,30]. Fully convolutional Siamese networks have also been proposed to produce dense and efficient sliding-window embeddings, with notable success on visual object tracking tasks [9,2,10]. Another popular technique for contrastive learning is triplet networks [12]. We also use Siamese networks to learn embeddings; however, we propose a process to learn embeddings that are invariant to spurious asymmetries, while being sensitive to pathology-inducing ones.

## 3  Method

### 3.1  Problem Setup

Given a PXR, denoted as $I$, we aim to detect sites of anterior pelvic fractures. Following the widely adopted approach by CAD methods [20,33,34], our model produces image-level binary classifications of fracture and heatmaps as fracture localization. Using heatmaps to represent localization (instead of bounding box or segmentation) stems from the inherent ambiguity in the definition of instance and boundary of pathological abnormalities in medical images. For instance, a fracture can be comminuted, *i.e.* bone breaking into multiple pieces, resulting in ambiguity in defining the number of fractures. Our model takes a cost-effective and flexible annotation format, a point at the center of each fracture site, allowing ambiguous fracture conditions to be flexibly represented as one point or multiple points. We dilate the annotation points by an empirically-defined radius (2 cm in our experiment) to produce a mask for the PXR, which is the training target of our method, denoted as $M$. In this way, we execute heatmap regression, similar to landmark detection [36], except for center-points of abnormalities with ambiguous extents.
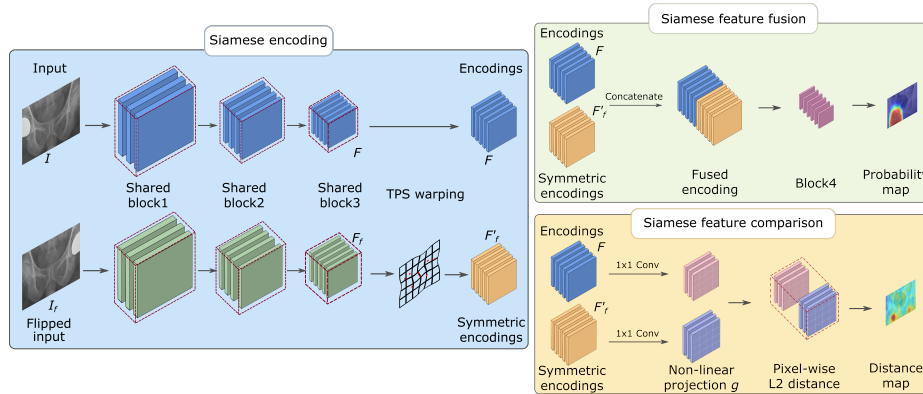
**Fig. 3.** System overview of the proposed AASN. The Siamese encoding module takes two pre-processed ROIs as input and encodes them using dense blocks with shared weights. After warping and alignment, the encoded feature maps are further processed by a Siamese feature fusion module and a Siamese contrastive learning module to produce a fracture probability map and a feature distance map, respectively.

### 3.2   Anatomy-Grounded Symmetric Warping

Given the input PXR image, our method first produces ROI of the anterior pelvis and anatomically-grounded warp to reveal the bilateral symmetry of the anatomy. The steps of ROI and warp generation are illustrated in Figure 2. First, a powerful graph-based landmark detection [19] is applied to detect 16 skeletal landmarks, including 7 pairs of bilateral symmetric landmarks and 2 points on pubic symphysis. From the landmarks, a line of bilateral symmetry is regressed, and the image is flipped with respect to it. Since we focus on detecting anterior pelvic fractures, where the dangers of massive bleeding is high and fractures are hard to detect, we extract ROIs of the anterior pelvis from the two images as a bounding box of landmarks on the pubis and ischium, which are referred as $I$ and $I_f$. A pixel-to-pixel warp from $I_f$ to $I$ is generated from the corresponding landmarks in $I_f$ and $I$ using the classic TPS warp [3], denoted as $T$. Note, the warp $T$ is not directly used to align the images. Instead, it is used in our Siamese network via a spatial transformer layer to align the features.

### 3.3   Anatomy-Aware Siamese Network

The architecture of AASN is shown in Figure 3. AASN contains a fully convolutional Siamese network with a DenseNet-121 [13] backbone. The dense blocks are split into two parts, an encoding part and a decoding part. It is worth noting that AASN allows the backbone network to be split flexibly at any block. For our application, we split at a middle level after the 3rd dense block, where the features are deep enough to encode the local skeletal pattern, but has not been pooled too heavily so that the textual information of small fractures is lost.
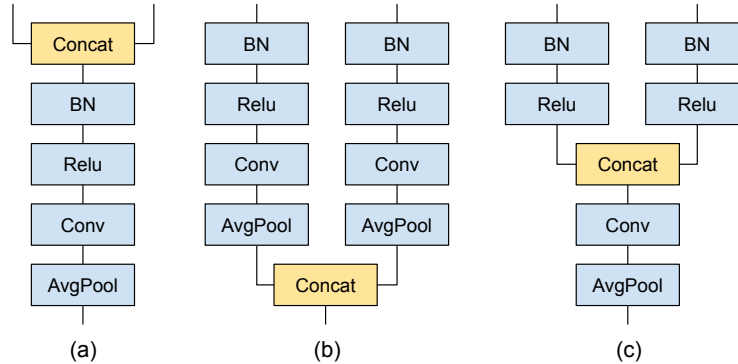
**Fig. 4.** Transition layer modification options for feature map fusion. (a) Feature map fusion before transition. (b) Feature map fusion after transition. (c) Feature map fusion inside transition

The encoding layers follow a Siamese structure, with two streams of weight-shared encoding layers taking the two images $I$ and $I_f$ as inputs. The encoder outputs, denoted as $F$ and $F_f$, provide feature representations of the original image and the flipped image, respectively. The spatial alignment transform $T$ is applied on $F_f$, resulting in $F'_f$, making corresponding pixels in $F$ and $F'_f$ represent corresponding anatomies. The two aligned feature maps are then fused and decoded to produce a fracture probability map, denoted as $Y$. Details of feature map fusion and decoding will be described in Sec. 3.4. We produce the probability heatmap as fracture detection result to alert the clinician the presence of a fracture and also to guide his or her attention (as shown in Figure 6). Since pelvis fractures can be very difficult to detect, even when there is a known fracture, this localization is a key feature over-and-above image-level predictions.

The model is trained using two losses. The first loss is the pixel-wise BCE between the predicted heatmap $Y$ and the ground truth $M$, denoted as $L_b$. The second loss is the pixel-wise contrastive loss between the two feature maps, $F$ and $F'_f$, denoted as $L_c$. Details of the contrastive loss will be discussed in Sec. 3.5. The total loss can be written as

$$L = L_b + \lambda L_c, \tag{1}$$

where $\lambda$ is a weight balancing the two losses.

### 3.4 Siamese Feature Fusion

The purpose of encoding the flipped image is to provide a reference of the symmetric counterpart, $F_f$, which can be incorporated with the feature $F$ to facilitate fracture detection. To provide a meaningful reference, $F_f$ needs to be spatially aligned with $F$, so that features with the same index/coordinate in the two feature maps encode the same, but symmetric, anatomies of the patient.

Previous methods have aligned the bilateral images $I$ and $I_f$ directly before encoding [24]. However, when large imaging angle and patient pose variations are present, image alignment is prone to introducing artifacts, which can increase the difficulty of fracture detection. Therefore, instead of aligning the bilateral images directly, we apply a spatial transformer layer on the feature map $F_f$ to align it with $F$, resulting in $F'_f$. The aligned feature maps $F$ and $F'_f$ are fused to produce a bilaterally combined feature map, where every feature vector encodes the visual patterns from symmetrical anatomies. This allows the decoder to directly incorporate symmetry analysis into fracture detection.

We fuse the feature maps by concatenation. Implementation of the concatenation involves modification to the transition module between the dense blocks, where multiple options exist, including concatenation before, after, or inside the transition module (as shown in Figure 4). A transition module in DenseNet consists of sequential BatchNorm, ReLU, Conv and AvgPool operations. We perform the concatenation inside the transition module after the ReLU layer, because it causes minimal structural changes to the DenseNet model. Specifically, the only layer affected in the DenseNet is the $1 \times 1$ Conv layer after concatenation, whose input channels are doubled. All other layers remain the same, allowing us to leverage the ImageNet pre-trained weights.

### 3.5   Siamese Contrastive Learning

While the above feature fusion provides a principled way to perform symmetric analysis, further advancements can be made. We are motivated by a key insight that image asymmetry can be caused by pathological abnormalities, *i.e.* fracture, or spurious non-pathological factors, *e.g.* soft tissue shadows, bowel gas patterns, clothing and foreign bodies. These non-pathological factors can be visually confusing, causing false positives. We aim to optimize the deep features to be more salient to the semantically pathological asymmetries, while mitigating the impact of distracting non-pathological asymmetries. To this end, our model employs a new constrastive learning component to minimize the pixel-wise distance between $F$ and $F'_f$ in areas without fracture, making the features insensitive to non-semantic asymmetries and thus less prone to false positives. On the other hand, our contrastive learning component encourages larger distance between $F$ and $F'_f$ in areas with fractures, making the features more sensitive to semantic asymmetries.

The above idea is implemented using pixel-wise margin loss between $F$ and $F'_f$ after a non-linear projection $g$:

$$L_c = \sum_{\boldsymbol{x}} \begin{cases} \|g(F(\boldsymbol{x})) - g(F'_f(\boldsymbol{x}))\|^2 & \text{if } \boldsymbol{x} \notin \hat{M} \\ \max(0, m - \|g(F(\boldsymbol{x})) - g(F'_f(\boldsymbol{x}))\|^2) & \text{if } \boldsymbol{x} \in \hat{M} \end{cases}, \qquad (2)$$

where $\boldsymbol{x}$ denotes the pixel coordinate, $\hat{M}$ denotes the mask indicating areas affected by fractures, and $m$ is a margin governing the dissimilarity of semantic asymmetries. The mask $\hat{M}$ is calculated as $\hat{M} = M \cup T \circ M_f$, where $T \circ M_f$ is flipped and warped $M$.

We employ a non-linear projection $g$ to transform the feature before calculating the distance, which improves the quality of the learned feature $F$, $F_f'$. In our experiment, the non-linear projection consists of a linear layer followed by BatchNorm and ReLU. We posit that directly performing contrast learning on features used for fracture detection could induce information loss and limit the modeling power. For example, bone curvature asymmetries in X-ray images are often non-pathological (e.g., caused by pose). However, they also provide visual cues to detect certain types of fractures. Using the non-linear projection, such useful information can be excluded from the contrastive learning so that they are preserved in the feature for the downstream fracture detection task.

While the margin loss has been adopted for CAD in a previous method [22], it was employed as a metric learning tool to learn a distance metric that directly represent the image asymmetry. We stress that our targeted CAD is more complex and clinically relevant, where image asymmetry can be semantically non-pathological (caused by pose, imaging condition and etc.) but we are only interested in detecting the pathological (fracture-caused) asymmetries. We employ the margin loss in our contrastive learning component to learn features with optimal properties. For this purpose, extra measures are taken in our method, including 1) conducting multi-task training with the margin loss calculated on a middle level feature, and 2) employing a non-linear projection head to transform the feature before calculating the margin loss.

## 4    Experiments

We demonstrate that our proposed AASN can significantly improve the performance in pelvic fracture detection by exploiting the semantic symmetry of anatomies. We focus on detecting fractures on the anterior pelvis including pubis and ischium, an anatomically symmetric region with high rate of diagnostic errors and life-threatening complications in the clinical practice.

### 4.1    Experimental Settings

**Dataset:**  We evaluate AASN on a real-world clinical dataset collected from the Picture Archiving and Communication System (PACS) of a hospital's trauma emergency department. The images have a large variation in the imaging conditions, including viewing angle, patient pose and foreign bodies shown in the images. Fracture sites in these images are labeled by experienced clinicians, combining multiple sources of information for confirmation, including clinical records and computed tomography scans. The annotations are provided in the form of points, due to inherent ambiguity in defining fracture as object. In total, there are 2 359 PXRs, and 759 of them have at least one anterior pelvic fracture site. All our experiments are conducted with five-fold cross-validation with a 70%/10%/20% training, validation, and testing split, respectively.

**Implementation Details:**  The ROIs of the anterior pelvis are resized to $256 \times 512$ and stacked to a 3-channel pseudo-color image. We produce the supervision mask for the heatmap prediction branch by dilating the annotation points
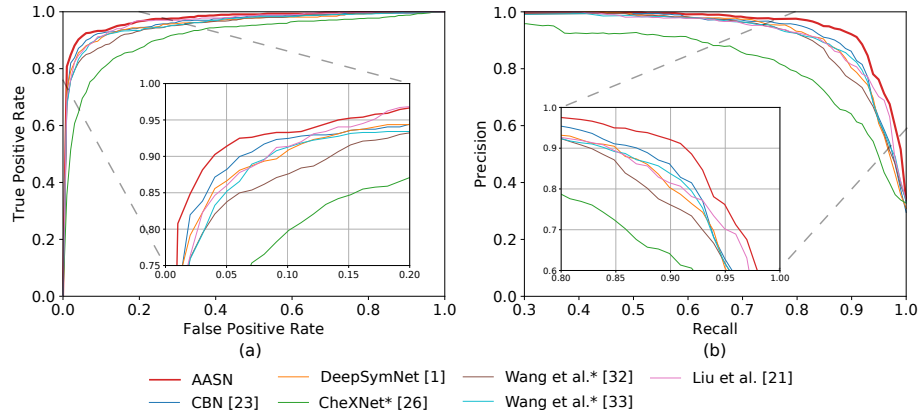
**Fig. 5.** Comparison of ROC curve and PR curve to the baselines. (a) is the ROC curve and (b) is the PR curve. *Methods trained using image-level labels.

to circle masks with a radius of 50 (about 2 cm). We implement all models using PyTorch [27]. Severe over-fitting is observed when training the networks from scratch, so we initialize them with ImageNet pre-trained weights. We emperically select DenseNet-121 as the backbone which yields the best performance comparing to other ResNet and DenseNet settings. All models are optimized by Adam [16] with a learning rate of $10^{-5}$. For the pixel-wise contrastive loss, we use the hyperparameter $m = 0.5$ as the margin, and $\lambda = 0.5$ to balance the total loss.

**Evaluation Metrics:** We first assess the model's performance as an image-level classsifier, which is a widely adopted evaluation approach for CAD systems [20,33,34]. The image-level abnormality reporting is of utmost importance in clinical workflow because it directly affects the clinical decision. We take the maximum value of the output heatmap as the classification output, and use Area under ROC Curve (AUC) and Average Precision (AP) to evaluate the classification performance.

We also evaluate the model's fracture localization performance. Since our model produces heatmaps as fracture localization, standard object detection metrics do not apply. A modified free-response ROC (FROC) is reported to measure localization performance. Specifically, unlike FROC, where object recall is reported with the number of false positives per image, we report fracture recall with the ratio of false positive area per image. A fracture is considered recalled if the heatmap activation value at its location is above the threshold. Areas with >2 cm away from all fracture annotation points are considered negative, on which the false positive ratio is calculated. Areas within 2 cm from any annotation point is considered as ambiguous extents of the fracture. Since both positive and negative responses in these ambiguous areas are clinically acceptable, they are excluded from the modified FROC calculation.

**Table 1.** Fracture classification and localization performance comparison with state-of-the-art models. Classifier AUC and AP are reported for classification performance. Fracture recalls at given false positive ratio are reported for localization performance. *Methods trained using image-level labels. Localization performance are not evaluated on these methods.

| Method | Classification | | Localization | |
|--------|------|------|----------------------|------------------------|
|        | AUC  | AP   | $\text{Recall}_{\text{FP}=1\%}$ | $\text{Recall}_{\text{FP}=10\%}$ |
| CheXNet* [28] | 93.42% | 86.33% | - | - |
| Wang *et al.*\* [34] | 95.43% | 93.31% | - | - |
| Wang *et al.*\* [35] | 96.06% | 93.90% | - | - |
| Liu *et al.* [22] | 96.84% | 94.29% | 2.78% | 24.19% |
| DeepSymNet [1] | 96.29% | 94.45% | 69.66% | 90.07% |
| CBN [24] | 97.00% | 94.92% | 73.93% | 90.90% |
| AASN | **97.71%** | **96.50%** | **77.87%** | **92.71%** |

**Compared Methods:**  We first compare AASN with three state-of-the-art CAD methods, *i.e.*, ChexNet [28], Wang *et al.* [34], and Wang *et al.* [35], all using image-level labels for training. They classify abnormality at image-level, and output heatmaps for localization visualization. ChexNet [28] employs a global average pooling followed by a fully connected layer to produce the final prediction. Wang *et al.* [34] uses Log-Sum-Exp (LSE) pooling. Wang *et al.* [35] employs a two-stage classification mechanism, and reports the state-of-the-art performance on hip/pelvic fracture classification.

We also compare with three methods modeling symmetry for CAD, *i.e.*, Liu *et al.* [22], CBN [24] and DeepSymNet [1]. All three methods perform alignment on the flipped image. Liu *et al.* [22] performs metric learning to learn a distance metric between symmetric body parts and uses it directly as an indicator of abnormalities. DeepSymNet [1] and CBN [24] fuse the Siamese encodings for abnormality detection, using subtraction and concatenation with gating, respectively. All evaluated methods use DenseNet-121 backbone, trained using the same experiment setting and tested with five-fold cross validation.

### 4.2   Classification Performance

Evaluation metrics of fracture classification performance are summarized in Table 1. ROC and PR curves are shown in Figure 5. The methods trained using only image-level labels result in overall lower performance than methods trained using fracture sites annotations. AASN outperforms all other methods, including the ones using symmetry and fracture site annotations, with substantial margins in all evaluation metrics. The improvements are also reflected in the ROC and PR curves Figure 5. Specifically, comparing to the 2nd highest values among all methods, AASN improves AUC and AP by 0.71% and 1.58%, from 97.00% and 94.92% to 97.71% and 96.50%, respectively. We stress that in this high AUC and

AP range (*i.e.* above 95%), the improvements brought by AASN are significant. For instance, when recall is increased from 95% to 96%, the number of missed fractures are reduced by 20%.

Figure 6 provides visualizations of fracture heatmaps produced using different methods. Non-displaced fractures that do not cause bone structures to be largely disrupted are visually ambiguous and often missed by the vanilla DenseNet-121 without considering symmetry. Comparison between the fracture site and its symmetric bone reveals that the suspicious pattern only occurs on one side and is likely to be fracture. This intuition is in line with the results, *i.e.*, by incorporating symmetric features, some of the ambiguous fractures can be detected. By employing the feature comparison module, AASN is able to detect more fracture, hypothetically owing to the better feature characteristics learned via feature comparison.

### 4.3   Localization Performance

We also evaluate AASN's fracture localization performance. The three symmetry modeling baselines and our four ablation study methods are also evaluated for comparison. As summarized in Table 1, AASN achieves the best fracture site recall among all evaluated methods, resulting in $\text{Recall}_{\text{FP}=1\%}$=77.87% and $\text{Recall}_{\text{FP}=10\%}$=92.71%, respectively. It outperforms baseline methods by substantial margins.

Among the baseline methods, directly using learned distance metric as an indicator of fracture (Liu *et al.* [22]) results in the lowest localization performance, because the image asymmetry indicated by distance metric can be caused by other non-pathological factors than fractures. The comparison justifies the importance of our proposed contrastive learning component, which *exploits image asymmetry to optimize deep feature for downstream fracture detection*, instead of directly using it as a fracture indicator. CBN [24] achieves the best performance among the three baselines, hypothetically owing to the Siamese feature fusion. With our feature alignment and contrastive learning components, AASN significantly improves fracture site $\text{Recall}_{\text{FP}=1\%}$ over CBN [24] by 3.94%.

### 4.4   Ablation Study

We conduct ablation study of AASN to analyze the contributions of its novel components, summarized in Table 2. The components include: 1) Symmetric feature fusion (referred to as *FF*), 2) Feature alignment (referred to as *FA*) and 3) Feature contrastive learning (referred to as *CL*). We add these components individually to the Vanilla DenseNet-121 to analyze their effects. We also analyze the effect of the non-linear projection head $g$ by evaluating a variant of constrastive learning without it.

**Symmetric Feature Fusion:** The effect of feature fusion is reflected in the comparisons: baseline vs. *FF* and baseline vs. *FF-FA*. Both *FF* and *FF-FA* employ symmetric feature fusion and are able to outperform *Vanilla*, although by a
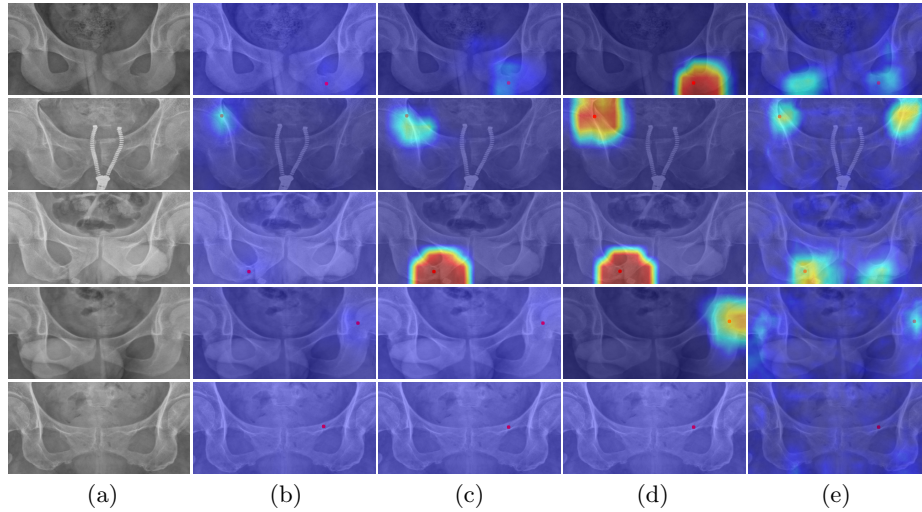
**Fig. 6.** Prediction results for different models. (a) pubis ROI in the PXR. Fracture probability heatmaps produced by (b) Vanilla DenseNet-121 [14], (c) CBN [24] and (d) AASN. (e) the distance map between siamese feature in AASN. The last row shows an example of failed cases.

different margin due to the different alignment methods used. In particular, *FF-FA* significantly improves the $\text{Recall}_{\text{FP}=1\%}$ by 5.89%. These improvements are hypothetically owing to the incorporation of the visual patterns from symmetric body parts, which provides reference for differentiating visually ambiguous fractures.

**Feature Alignment:** The effect of feature warping and alignment is reflected in the comparisons: *FF* vs. *FF-FA* and *FF-CL* vs. *FF-FA-CL*. The ablation study shows that, by using the feature warping and alignment, the performances of both *FF* and *FF-CL* are both significantly improved. In particular, the $\text{Recall}_{\text{FP}=1\%}$ are improved by 3.46% and 1.60% in *FF-FA* and *FF-FA-CL*, respectively. It's also demonstrated that the contributions of feature warping and alignment are consistent with and without Siamese feature comparison. We posit that the performance improvements are owing to the preservation of the original image pattern by performing warping and alignment at the feature level.

**Contrastive Learning:** The effect of Siamese feature comparison is reflected in the comparisons: *FF* vs. *FF-CL* and *FF-FA* vs. *FF-FA-CL*. The ablation study shows measurable contribution of the Siamese feature comparison module. By using Siamese feature fusion, *FF* and *FF-FA* already show significant improvements comparing to the baseline. By adding Siamese feature comparison to *FF* and *FF-FA*, $\text{Recall}_{\text{FP}=1\%}$ are improved by 3.05% and 1.19%, respectively. The improvements are in line with our motivation and hypothesis that by maximizing/minimizing Siamese feature distances on areas with/without fractures, the

**Table 2.** Ablation study of AASN. The baseline model is vanilla DenseNet121 trained without the symmetry modeling components. "FF" indicates using feature fusion. "FA" indicates using feature alignment (otherwise image alignment is used). "CL" indicates using contrastive learning. "no. proj." indicates that the contrastive learning is performed without the non-linear projection head.

| FF | FA | CL | AUC | AP | $\text{Recall}_{\text{FP}=1\%}$ | $\text{Recall}_{\text{FP}=10\%}$ |
|----|----|----|-----|----|-----|-----|
| | | | 96.52% | 94.52% | 70.79% | 89.46% |
| ✓ | | | 96.93% (+0.41%) | 94.77% (+0.25%) | 73.22% (+2.43%) | 89.93% (+0.47%) |
| ✓ | ✓ | | 97.20% (+0.68%) | 95.68% (+1.16%) | 76.68% (+5.89%) | 91.51% (+2.05%) |
| ✓ | | ✓ | 97.46% (+0.94%) | 95.36% (+0.84%) | 76.27% (+5.48%) | 91.09% (+1.63%) |
| ✓ | ✓ | ✓ no proj. | 97.31% (+0.79%) | 96.15% (+1.63%) | 77.26% (+6.47%) | 92.70% (+3.24%) |
| ✓ | ✓ | ✓ | **97.71%** (+1.19%) | **96.50%** (+1.98%) | **77.87%** (+7.08%) | **92.71%** (+3.25%) |

network can learn features that are more sensitive to fractures and less sensitive to other distracting factors. Comparing to the AASN directly performing contrastive learning on the symmetric feature (*no. proj.*), employing the non-linear projection head leads further improves the $\text{Recall}_{\text{FP}=1\%}$ by 0.61%.

## 5    Conclusion

In this paper, we systematically and thoroughly study exploiting the anatomical symmetry prior knowledge to facilitate CAD, in particular anterior pelvic fracture detection in PXR. We introduce a deep neural network technique, termed Anatomical-Aware Siamese Network, to incorporate semantic symmetry analysis into abnormality (*i.e.* fracture) detection. Through comprehensive ablation study, we demonstrate that: 1) Employing symmetric feature fusion can effectively exploit symmetrical information to facilitate fracture detection. 2) Performing spatial alignment at the feature level for symmetric feature fusion leads to substantial performance gain. 3) Using contrastive learning, the Siamese encoder is able to learn more sensible embedding, leading to further performance improvement. By comparing with the state-of-the-art methods, including latest ones modeling symmetry, we demonstrate the AASN is by far the most effective method exploiting symmetry and reports substantially improved performances on both classification and localization tasks.

# References

1. Barman, A., Inam, M.E., Lee, S., Savitz, S., Sheth, S., Giancardo, L.: Determining ischemic stroke from ct-angiography imaging using symmetry-sensitive convolutional networks. In: 2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019). pp. 1873–1877 (April 2019). https://doi.org/10.1109/ISBI.2019.8759475
2. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.: Fully-convolutional siamese networks for object tracking. In: European conference on computer vision. pp. 850–865. Springer (2016)
3. Bookstein, F.L.: Principal warps: thin-plate splines and the decomposition of deformations. IEEE Transactions on Pattern Analysis and Machine Intelligence **11**(6), 567–585 (June 1989). https://doi.org/10.1109/34.24792
4. Bustos, A., Pertusa, A., Salinas, J.M., de la Iglesia-Vay, M.: Padchest: A large chest x-ray image dataset with multi-label annotated reports (2019)
5. Chen, H., Miao, S., Xu, D., Hager, G.D., Harrison, A.P.: Deep hierarchical multi-label classification of chest x-ray images. In: Cardoso, M.J., Feragen, A., Glocker, B., Konukoglu, E., Oguz, I., Unal, G., Vercauteren, T. (eds.) Proceedings of The 2nd International Conference on Medical Imaging with Deep Learning. Proceedings of Machine Learning Research, vol. 102, pp. 109–120. PMLR, London, United Kingdom (08–10 Jul 2019), `http://proceedings.mlr.press/v102/chen19a.html`
6. Cheng, C.T., Ho, T.Y., Lee, T.Y., Chang, C.C., Chou, C.C., Chen, C.C., Chung, I.F., Liao, C.H.: Application of a deep learning algorithm for detection and visualization of hip fractures on plain pelvic radiographs. European radiology pp. 1–9 (2019)
7. Clohisy, J.C., Carlisle, J.C., Beaulé, P.E., Kim, Y.J., Trousdale, R.T., Sierra, R.J., Leunig, M., Schoenecker, P.L., Millis, M.B.: A systematic approach to the plain radiographic evaluation of the young adult hip. The Journal of Bone and Joint Surgery. American volume. **90**(Suppl 4),  47 (2008)
8. Gale, W., Oakden-Rayner, L., Carneiro, G., Bradley, A.P., Palmer, L.J.: Detecting hip fractures with radiologist-level performance using deep neural networks. CoRR **abs/1711.06504** (2017), `http://arxiv.org/abs/1711.06504`
9. Guo, Q., Feng, W., Zhou, C., Huang, R., Wan, L., Wang, S.: Learning dynamic siamese network for visual object tracking. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1763–1771 (2017)
10. Guo, Q., Feng, W., Zhou, C., Huang, R., Wan, L., Wang, S.: Learning dynamic siamese network for visual object tracking. In: The IEEE International Conference on Computer Vision (ICCV) (Oct 2017)
11. Hadsell, R., Chopra, S., LeCun, Y.: Dimensionality reduction by learning an invariant mapping. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06). vol. 2, pp. 1735–1742 (June 2006). https://doi.org/10.1109/CVPR.2006.100
12. Hoffer, E., Ailon, N.: Deep metric learning using triplet network. In: Feragen, A., Pelillo, M., Loog, M. (eds.) Similarity-Based Pattern Recognition. pp. 84–92. Springer International Publishing, Cham (2015)
13. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.: Densely connected convolutional networks. arxiv website. arxiv. org/abs/1608.06993. Published August **24** (2016)
14. Huang, G., Liu, Z., Weinberger, K.Q.: Densely connected convolutional networks. CoRR **abs/1608.06993** (2016), `http://arxiv.org/abs/1608.06993`

15. Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R.L., Shpanskaya, K., Seekins, J., Mong, D.A., Halabi, S.S., Sandberg, J.K., Jones, R., Larson, D.B., Langlotz, C.P., Patel, B.N., Lungren, M.P., Ng, A.Y.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. CoRR **abs/1901.07031** (2019)

16. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)

17. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In: ICML deep learning workshop. vol. 2 (2015)

18. Konukoglu, E., Glocker, B., Criminisi, A., Pohl, K.M.: Wesd–weighted spectral distance for measuring shape dissimilarity. IEEE transactions on pattern analysis and machine intelligence **35**(9), 2284–2297 (2012)

19. Li, W., Lu, Y., Zheng, K., Liao, H., Lin, C., Luo, J., Cheng, C.T., Xiao, J., Lu, L., Kuo, C.F., et al.: Structured landmark detection via topology-adapting deep graph learning. arXiv preprint arXiv:2004.08190 (2020)

20. Li, Z., Wang, C., Han, M., Xue, Y., Wei, W., Li, L.J., Fei-Fei, L.: Thoracic Disease Identification and Localization with Limited Supervision. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8290–8299. IEEE, Salt Lake City, UT (Jun 2018). https://doi.org/10.1109/CVPR.2018.00865

21. Ling, H., Gao, J., Kar, A., Chen, W., Fidler, S.: Fast interactive object annotation with curve-gcn. CoRR **abs/1903.06874** (2019), `http://arxiv.org/abs/1903.06874`

22. Liu, C.F., Padhy, S., Ramachandran, S., Wang, V.X., Efimov, A., Bernal, A., Shi, L., Vaillant, M., Ratnanather, J.T., Faria, A.V., Caffo, B., Albert, M., Miller, M.I.: Using deep siamese neural networks for detection of brain asymmetries associated with alzheimer's disease and mild cognitive impairment. Magnetic Resonance Imaging (2019). https://doi.org/https://doi.org/10.1016/j.mri.2019.07.003, `http://www.sciencedirect.com/science/article/pii/S0730725X19300086`

23. Liu, S.X.: Symmetry and asymmetry analysis and its implications to computer-aided diagnosis: A review of the literature. Journal of biomedical informatics **42**(6), 1056–1064 (2009)

24. Liu, Y., Zhou, Z., Zhang, S., Luo, L., Zhang, Q., Zhang, F., Li, X., Wang, Y., Yu, Y.: From unilateral to bilateral learning: Detecting mammogram masses with contrasted bilateral network. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. pp. 477–485. Springer International Publishing, Cham (2019)

25. Lu, Y., Zheng, K., Li, W., Wang, Y., Harrison, A.P., Lin, C., Wang, S., Xiao, J., Lu, L., Kuo, C.F., et al.: Learning to segment anatomical structures accurately from one exemplar. arXiv preprint arXiv:2007.03052 (2020)

26. Melekhov, I., Kannala, J., Rahtu, E.: Siamese network features for image matching. In: 2016 23rd International Conference on Pattern Recognition (ICPR). pp. 378–383 (Dec 2016). https://doi.org/10.1109/ICPR.2016.7899663

27. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)

28. Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D.Y., Bagul, A., Langlotz, C., Shpanskaya, K.S., Lungren, M.P., Ng, A.Y.: Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. CoRR **abs/1711.05225** (2017), `http://arxiv.org/abs/1711.05225`

29. Simonovsky, M., Gutiérrez-Becker, B., Mateus, D., Navab, N., Komodakis, N.: A deep metric for multimodal registration. In: International conference on medical image computing and computer-assisted intervention. pp. 10–18. Springer (2016)
30. Sun, Y., Wang, X., Tang, X.: Deep learning face representation by joint identification-verification. CoRR **abs/1406.4773** (2014), `http://arxiv.org/abs/1406.4773`
31. Varior, R.R., Haloi, M., Wang, G.: Gated siamese convolutional neural network architecture for human re-identification. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) Computer Vision – ECCV 2016. pp. 791–808. Springer International Publishing, Cham (2016)
32. Wachinger, C., Golland, P., Kremen, W., Fischl, B., Reuter, M., Initiative, A.D.N., et al.: Brainprint: A discriminative characterization of brain morphology. NeuroImage **109**, 232–248 (2015)
33. Wang, H., Xia, Y.: Chestnet: A deep neural network for classification of thoracic diseases on chest radiography. arXiv preprint arXiv:1807.03058 (2018)
34. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
35. Wang, Y., Lu, L., Cheng, C.T., Jin, D., Harrison, A.P., Xiao, J., Liao, C.H., Miao, S.: Weakly supervised universal fracture detection in pelvic x-rays. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. pp. 459–467. Springer International Publishing, Cham (2019)
36. Xu, Z., Huo, Y., Park, J., Landman, B., Milkowski, A., Grbic, S., Zhou, S.: Less is more: Simultaneous view classification and landmark detection for abdominal ultrasound images. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 711–719. Springer (2018)
37. Zagoruyko, S., Komodakis, N.: Learning to compare image patches via convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4353–4361 (2015)
38. Zhou, B., Khosla, A., Lapedriza, À., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. CoRR **abs/1512.04150** (2015), `http://arxiv.org/abs/1512.04150`