

# 000 001 002 003 004 005 POSE-GUIDED PROPORTIONAL STRUCTURE ASSESS- 006 MENT FOR 2D ILLUSTRATED CHARACTERS 007 008

009 **Yu-Jyuan Lin<sup>1</sup>, Chian-Ru Lin<sup>2</sup>, Si-Xian Hsu<sup>2</sup>**  
 010  
011  
012

013 <sup>1</sup>Department of Information Management, National Taiwan University  
 014  
015

016 <sup>2</sup>Department of Electrical Engineering, National Taiwan University  
 017 {b13705012, b13901007, b13901009}@ntu.edu.tw  
 018  
019  
020  
021  
022  
023

## ABSTRACT

013 Structural errors in human anatomy, particularly proportion and composition, are  
 014 common challenges for beginners in 2D character illustration. Current research  
 015 primarily focuses on content generation rather than structural correction and diag-  
 016 nosis. This study proposes a novel evaluation framework for detecting abnormal  
 017 proportions in illustrated figures. Our method combines a Spatial Transformer  
 018 Network for geometric alignment with a Siamese CNN for proportion-aware sim-  
 019 ilarity learning, and uses Grad-CAM to visualize matching results as guidance  
 020 for beginner illustrators. This framework provides an interpretable and objective  
 021 approach for assessing proportional errors in 2D illustrations, offering practical  
 022 guidance for improving structural accuracy in manga character drawing.  
 023

## 1 INTRODUCTION

024 For beginners in anime illustration, mastering human anatomy and skeletal structure poses a signif-  
 025 icant challenge. Novices often fall into the trap of local depiction, focusing excessively on details  
 026 while neglecting the overall spatial arrangement. This frequently results in unbalanced skeletal  
 027 structures or collapsed perspectives in the final artwork. Unlike realistic drawing, which adheres to  
 028 fixed anatomical standards, 2D anime creation allows for a certain degree of stylistic variation. This  
 029 flexibility makes the criteria for incorrect anatomy ambiguous and difficult to define.

030 In the field of computer vision, current research on 2D anime characters focuses primarily on col-  
 031 orization, line art extraction, or generation, with less exploration into drawing tutorials or skeletal  
 032 correction. While mature technologies for pose extraction and rehabilitation action assessment exist  
 033 for natural images, relevant literature in the anime domain remains scarce. To address this, we sim-  
 034 ulate proportional abnormalities and diverse compositions using scaling, translation, and warping  
 035 techniques to model common learner errors.

036 To address this gap, we propose a proportionality-aware skeleton analysis framework based on an  
 037 Siamese convolutional neural network with the spatial transformer network (STN) [2](Jaderberg  
 038 et al., 2015) as an alignment module coupled with gradient-weighted Class Activation Mapping  
 039 (Grad-CAM) [4](Selvaraju et al., 2017)visualization. During inference, a user-provided skeleton  
 040 with potential proportional errors is automatically matched against a database of correctly propor-  
 041 tioned reference skeletons. The input skeleton and its most similar reference are then passed through  
 042 the two branches of the Siamese CNN to compute their cosine similarity, after which Grad-CAM is  
 043 applied to identify the regions most responsible for the similarity score. Consequently, the result-  
 044 ing heatmap highlights the proportionally correct regions of the input skeleton while suppressing  
 045 activation in distorted areas, offering intuitive and interpretable feedback for learners.  
 046

## 2 RELATED WORK

### 2.1 POSE-GUIDED MATCHING

047 We recognize that our task—detecting abnormal proportions in illustrated figures—shares the same  
 048 high-level idea as pose assessment in sports and rehabilitation, a well-developed research field. The  
 049  
050  
051  
052  
053

054 following section outlines related work that is applicable to our scenario with certain adjustments  
 055 and modifications.

056 One relevant paper focuses on aiding patients in rehabilitation exercises [3](Qiu et al., 2022). The  
 057 rehabilitation exercise of interest is the eight-section brocade, a traditional Chinese practice com-  
 058 posed of only eight key postures. Their training data consist of frames captured from videos of  
 059 patients and a video of a specialist performing the eight-section brocade. Only frames around the  
 060 key posture points are extracted, so the dataset can be categorized into eight classes. Images of the  
 061 specialist are treated as “standard,” while images of patients are treated as “learner.”

062 The paper introduces a new idea for pose matching, namely pose-guided matching, which aims  
 063 to provide objective and accurate scores, feedback, and guidance to patients when their poses are  
 064 compared with the standard poses. More specifically, the authors propose a pair-based Siamese  
 065 Convolutional Neural Network (SCNN), abbreviated ST-AMCNN, to realize pose-guided matching.  
 066 They simplify multi-stage pose matching by merging the alignment and matching modules into a  
 067 one-stage task, such that only one loss function is required, reducing computational complexity.  
 068 Building upon the Spatial Transformer Networks (STN) used as the alignment module, they propose  
 069 a new Attention-based Multi-Scale Convolution (AMC) to match different posture parts (i.e., multi-  
 070 scale). Finally, Gradient-weighted Class Activation Mapping (Grad-CAM) is adopted to visualize  
 071 the matching results for the learner.

072 This paper provides a workflow that we consider plausible to transplant to our scenario. However,  
 073 several challenges arise when adapting this standard operating procedure from a rehabilitation as-  
 074 sistant for patients to an illustration assistant for art beginners. First, it is impossible to obtain many  
 075 poorly illustrated figures paired with corrected versions drawn by experienced artists, so we must  
 076 create a method for generating (standard, learner) counterparts in our scenario. Second, we do not  
 077 aim to fit all standing-pose skeletons to a single standard pose, as this would violate the artist’s in-  
 078 tention. Instead, we aim to preserve most aspects of the original pose (e.g., facing direction, stance  
 079 width) while correcting only proportion errors, making our task significantly more challenging than  
 080 in the eight-section brocade case.

## 082 2.2 POSE ESTIMATION

084 We also examine a paper on illustrated character skeleton generation and dataset construction  
 085 [1](Chen & Zwicker, 2022). This work expands existing datasets and proposes a transfer-learning  
 086 model to address the scarcity of pose-estimation data for illustrations. In addition to achieving  
 087 high-accuracy skeleton extraction for illustrated figures, the authors further apply their state-of-the-  
 088 art character pose estimator to the novel task of pose-guided illustration retrieval. They construct  
 089 a retrieval system that searches for illustrated characters in similar poses by performing a simple  
 090 nearest-neighbor search of euclidean distances over normalized keypoints. The original purpose  
 091 of this system is to serve as a practical reference tool for artists, who commonly rely on reference  
 092 drawings during the illustration process.

093 This paper provides us with a reliable skeleton generator and a large dataset equipped with a retrieval  
 094 mechanism, thereby enabling more effective construction of a dataset tailored to our task.

## 096 3 METHODOLOGY

### 099 3.1 OVERVIEW

101 We construct a dataset of manga character skeletons for training and evaluation. Because it is not  
 102 feasible to collect a sufficient number of poorly illustrated figures paired with corrected versions  
 103 drawn by experienced artists, we generate our own ground truth. In this setup, the skeleton of  
 104 a well-illustrated figure serves as the standard, and warped variants of this skeleton serve as the  
 105 learners. Additional details regarding dataset generation are provided in the experiment section.

106 Our model is trained on paired images using a Siamese CNN as the primary architecture. Each image  
 107 pair is processed in the way that only the learner skeleton passes through a Spatial Transformer  
 Network (STN) module while the standard skeleton directly to the Siamese CNN. The Siamese

108 CNN outputs a pair of corresponding embeddings, which are used to compute a single loss value.  
 109 This loss is then backpropagated to jointly update both the CNN parameters and the STN parameters.  
 110

111 For the overall model design, we adopt the ST-AMCNN structure described in [3](Qiu et al., 2022).  
 112 This framework employs STN for alignment, a Siamese CNN for matching, and Grad-CAM for  
 113 generating visualized correction guidance. We retain most of the original structure and the hyperpa-  
 114 rameters due to the conceptual similarity between our task and the pose-assessment task. Nonethe-  
 115 less, several components are modified to better accommodate our application scenario and dataset  
 116 characteristics

### 117 3.2 SPATIAL TRANSFORMER NETWORK (STN)

119 Before the matching stage, the two skeletons being compared must be aligned to the same position,  
 120 size, and orientation. This ensures that the Siamese CNN focuses on differences in posture rather  
 121 than irrelevant variations in scale, location, or rotation. The structure of the STN is shown in Fig. 1.

122 The STN consists of three main components: a localization network, a grid generator, and a grid  
 123 sampler. The localization network is a CNN that produces a parameter vector  $\theta$  Figure. 2. The  
 124 grid generator then uses  $\theta$  to construct a transformation matrix  $A_\theta$ (Eq. 1a 1b). This matrix maps  
 125 the target coordinates back to the learner’s coordinate space; given this mapping, the reverse trans-  
 126 formation can also be obtained. The training objective of the STN is therefore to warp the input  
 127 skeleton into a centered, upright position.

### 128 3.3 SIAMESE CNN

131 We implement a Siamese CNN architecture composed of two Attention-based Multi-Scale Convo-  
 132 lutional Networks (AMCNN). AMCNN contains a bundle of three convolution kernels with differ-  
 133 ent sizes, designed to capture features at multiple spatial scales. The extracted features are then  
 134 passed into an attention module, which resembles the Squeeze-and-Excitation (SE) block proposed  
 135 in SENet. This attention mechanism encourages the network to focus on posture-relevant features  
 136 while suppressing irrelevant ones.

### 137 3.4 GRAD-CAM

139 To inform the learner which specific parts of the skeleton require correction, we apply Grad-CAM  
 140 to generate a heatmap. Grad-CAM takes the output of the last convolutional layer and computes the  
 141 cosine similarity between the two outputs. This visualizes the comparison result of the two input  
 142 poses from the CNN’s perspective. In the heatmap, red regions indicate that the corresponding parts  
 143 of the learner skeleton are more similar to the standard skeleton.

### 144 3.5 LOSS FUNCTION

147 Here, we used two loss function method for identify the model. We implement Cosine Embedding  
 148 Loss and evaluate its performance. It is important to note that the loss function penalizes not only  
 149 mismatches between the two skeletons, but also misalignment of the input skeleton produced by the  
 150 STN.

151 To compute similarity between the two embeddings, we use a cosine embedding loss. Because  
 152 cosine similarity measures directional alignment, it is more sensitive to posture orientation than  
 153 Euclidean distance. The embedding loss is defined as:

$$155 \quad L(x_1, x_2, y) = \begin{cases} 1 - \cos(x_1, x_2), & y = 1, \\ 156 \quad \max(0, \cos(x_1, x_2) - m), & y = -1, \end{cases}$$

157 where  $x_1$  and  $x_2$  are the embeddings corresponding to the standard and learner skeletons,  $y = 1$   
 158 indicates that the pair belongs to the same class,  $y = -1$  indicates different classes, and  $m$  is the  
 159 margin.

161 To avoid penalty, the model must learn to produce highly similar embeddings for same-class pairs  
 162 and sufficiently different embeddings for different-class pairs.

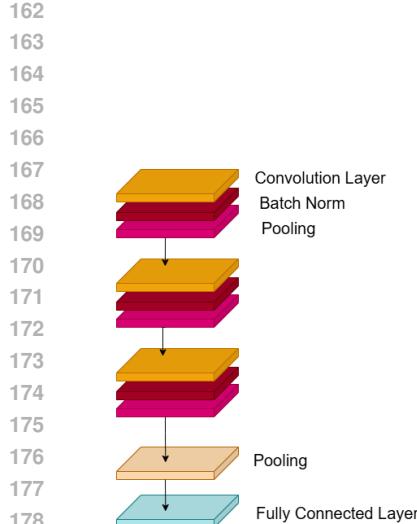


Figure 1: Localization Network of STN Module

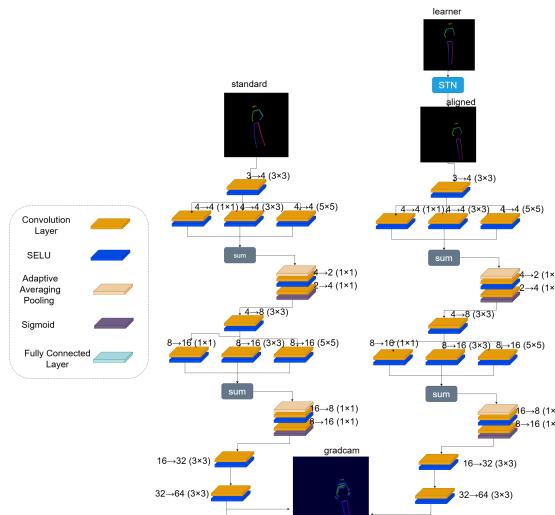


Figure 2: Architecture of our Training Model

$$\begin{bmatrix} x_i^s \\ y_i^s \\ 1 \end{bmatrix} = A_\theta \begin{bmatrix} x_i^t \\ y_i^t \\ 1 \end{bmatrix} = \begin{bmatrix} s_x \cos \theta & -s_x \sin \theta & t_x \\ s_y \sin \theta & s_y \cos \theta & t_y \end{bmatrix} \begin{bmatrix} x_i^t \\ y_i^t \\ 1 \end{bmatrix} \quad (1a)$$

$$s_x = 1 + \alpha s_x^{raw}, \quad s_y = 1 + \alpha s_y^{raw} \quad (1b)$$

**Equation 1a, 1b:** Affine transformation of STN.  $s_x, s_y$  are the scaling parameters,  $\alpha$  is the scaling factor,  $t_x, t_y$  are translation parameters,  $\theta$  is the rotation parameter.  $A_\theta$  is the affine matrix,  $(x_i^s, y_i^s)$  is the source coordinate,  $(x_i^t, y_i^t)$  is the target coordinate.

### 3.6 EVALUATION

To assess whether the model correctly identifies proportional errors, we construct test skeletons in which the disproportion is applied exclusively to one of four regions: the left half, right half, upper half, or lower half. For each test sample, we compute the cosine similarity score between the distorted skeleton and its corresponding standard skeleton, and then apply Grad-CAM to visualize which regions contribute most strongly to the score. Since the distorted region should not resemble the standard skeleton, an ideal model should produce weak or no activation on the distorted half, while highlighting structural correspondences on the undistorted half.

A prediction is therefore considered correct when the Grad-CAM activation on the undistorted region exceeds the activation on the distorted region, either in intensity or spatial extent. The final accuracy metric is computed as the proportion of test samples that satisfy this criterion.

Based on the proposed STN-Siamese CNN-Grad-CAM architecture, our experiments aim to evaluate whether the model can reliably identify skeletons with incorrect proportions, and whether spatial alignment via the STN further improves this capability. This section presents the construction of the dataset, the implementation details of our training setup, the comparative performance of models with and without the STN module, and the effects of varying key hyperparameters.

216 **4 EXPERIMENTS AND RESULTS**

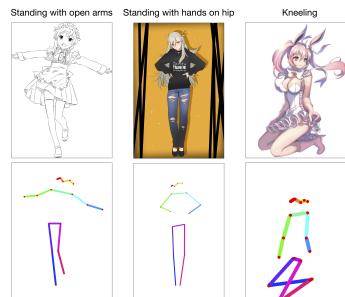
217 **4.1 DATASET**

218 We construct a dataset of manga-style skeletons specifically designed to ensure that the STN learns  
 219 meaningful geometric transformations while the subsequent Siamese CNN learns to discriminate  
 220 proportional correctness. Because the original dataset contains only a limited number of clean,  
 221 well-proportioned skeletons, we treat each of these as a standard reference and generate two com-  
 222plementary categories of training samples. The first category consists of eight proportion-preserving  
 223transformations that encourage the STN to learn alignment-related transformations without altering  
 224the underlying anatomy. The second category consists of nine proportion-changing distortions, en-  
 225abling the CNN to learn explicit differences between correct and incorrect proportions. Together,  
 226these two categories constitute approximately equal portions of the dataset, allowing the model to  
 227simultaneously acquire robust transformation invariance and proportion-awareness.

228 The training data are built from three canonical poses: standing with open arms, standing with  
 229 hands on hip, and kneeling(Fig. 3). We follow the pipeline described in Transfer Learning for Pose  
 230 Estimation of Illustrated Characters [1](Chen & Zwicker, 2022). For each pose, we first use the  
 231 provided pose-retrieval model to obtain candidate images whose pose similarity meets our criterion,  
 232 and then apply the released pose-estimation model to extract skeletons. To prevent confusion during  
 233 training, we manually filter out skeletons that deviate substantially from the intended pose cluster,  
 234 along with samples in which the character is facing away from the viewer. After filtering, we obtain  
 235 156 skeletons for the standing-with-open-arms pose, 141 for standing-with-hands-on-hip, and 120  
 236 for kneeling.  
 237

238 From each original skeleton, we generate a total of 17 transformed variants, resulting in 18 sam-  
 239 ples per original skeleton. Eight of these transformations preserve overall body proportions and  
 240 are designed to train the STN to compensate for global misalignments such as translation, rotation,  
 241 and scale variation. We apply two independent scaling augmentations to every original skeleton:  
 242 a random down-scaling to 80–100% of its original size and a random up-scaling to 100–120%. In  
 243 addition, each skeleton undergoes two independent rotation augmentations: a clockwise rotation  
 244 randomly sampled between 0° and 30°, and a counterclockwise rotation randomly sampled between  
 245 0° and 30°. To simulate positional shifts without causing the skeleton to exceed image boundaries,  
 246 we first down-scale the skeleton by 20% and then translate it by 60–90 pixels upward, downward,  
 247 leftward, or rightward. All proportion-preserving transformations maintain the relative ratios be-  
 248 tween body parts.

249 To expose the model to disproportionate skeletons, we apply 9 proportion-changing transformations  
 250 to each sample. These operations act on halves of the body and widen or narrow their width or  
 251 height asymmetrically. The scaling values are chosen so that the deformations are visually notice-  
 252 able while remaining anatomically plausible. In addition to asymmetric scaling, we incorporate a  
 253 swirl deformation in which the original image is warped by rotating pixels around the image center  
 254 with a maximum angle of roughly 30° that decays smoothly to zero at radius R. To preserve limb  
 255 straightness in the final representation, the swirl deformation is applied to the image first, after which  
 256 we run pose estimation to obtain the distorted skeleton.



257  
 258 **Figure 3: The sample of dataset**

270 4.2 IMPLEMENTATION DETAILS  
271

272 All models are trained on a combination of Google Colab L4, Kaggle P100, and Kaggle T4x2 GPUs.  
273 Unless otherwise mentioned, we use the Adam optimizer with 0.001 learning rate and batch size 32,  
274 for 100 epochs. The dataset is split into training, validation, and test sets in a 7:1:2 ratio.

275 4.3 ANALYSIS OF STN TRANSFORMATIONS  
276

277 To understand the influence of the STN on this task, we analyze the skeletons produced by the STN  
278 and the corresponding affine-transformation parameters. Both visual inspection and parameter anal-  
279 ysis indicate that the STN consistently learns to adjust translation and scale: skeletons that differ  
280 substantially in position or size are often shifted or rescaled toward a narrower range of configura-  
281 tions. In contrast, the rotation components of the learned affine matrices remain close to the identity  
282 transformation, even when the input has been rotated by up to 30°, suggesting that the STN makes  
283 limited use of rotational degrees of freedom under our training objective.

284 We further observe that the STN tends to translate skeletons toward the corners of the image, and in  
285 some cases, parts of the skeleton extend beyond the visible frame. Importantly, however, the Grad-  
286 CAM heatmaps produced from these outputs remain structurally aligned with the underlying skele-  
287 tons: high-activation regions consistently follow limb segments and joints regardless of whether the  
288 skeleton is centered or partially shifted. Consequently, the lack of explicit centering does not com-  
289 promise the interpretability or reliability of the heatmaps for proportion-related evaluation. Overall,  
290 the STN behaves as a coarse normalizer of size and position rather than a strict geometric aligner,  
291 and its irregular translations do not prevent the model from focusing on the structural cues relevant  
292 to proportionality.

293  
294 Table 1: The accuracy evaluation results of Section 4.4 and 4.5.

295 Class of distortion (half × scale)	296 STN	297 STN_bend	298 Deterministic	299 Deterministic_bend
297 Right × 0.5	15.79%	19.32%	21.43%	35.37%
298 Right × 1.3	6.25%	14.81%	20.48%	22.08%
299 Left × 0.7	19.00%	26.32%	34.29%	9.78%
300 Left × 1.5	10.26%	13.58%	15.91%	10.00%
301 Lower × 0.5	35.96%	9.20%	36.25%	16.51%
302 Lower × 1.3	23.16%	15.49%	52.75%	39.51%
303 Upper × 0.7	13.48%	16.67%	2.60%	11.11%
304 Upper × 1.5	9.64%	20.29%	3.41%	4.88%
<b>305 Aggregated accuracy by distortion direction</b>				
306 Right	22.04%	34.13%	41.91%	57.44%
307 Left	29.26%	39.90%	50.19%	19.78%
308 Lower	59.11%	24.59%	89.00%	56.02%
309 Upper	23.12%	36.96%	6.01%	15.99%
310 <b>Standard Deviation</b>	21.24%	6.59%	34.06%	22.49%
311 <b>Average</b>	16.69%	16.96%	23.39%	18.65%

312  
313 4.4 EFFECT OF STN-BASED ALIGNMENT  
314

315 To evaluate the impact of the STN on skeleton alignment and proportionality judgment, we compare  
316 the STN-based model with a deterministic alignment strategy. In the deterministic method, the  
317 central body axis—defined as the line segment connecting the midpoint of the shoulder keypoints  
318 to the midpoint of the pelvis keypoints—is transformed to a fixed position, orientation, and length  
319 for all skeletons. Because this alignment removes variations in translation, rotation, and global  
320 scale, only proportion-changing transformations are retained when constructing the baseline dataset,  
321 resulting in 10 variants per original skeleton. This baseline dataset bypasses the STN entirely and is  
322 fed directly into the Siamese CNN. The accuracy evaluation is shown in Table 1.

323 Although the axis-aligned baseline achieves a higher mean accuracy, its performance exhibits ex-  
324 tremely large variance across distortion types. The accuracy difference between detecting lower-

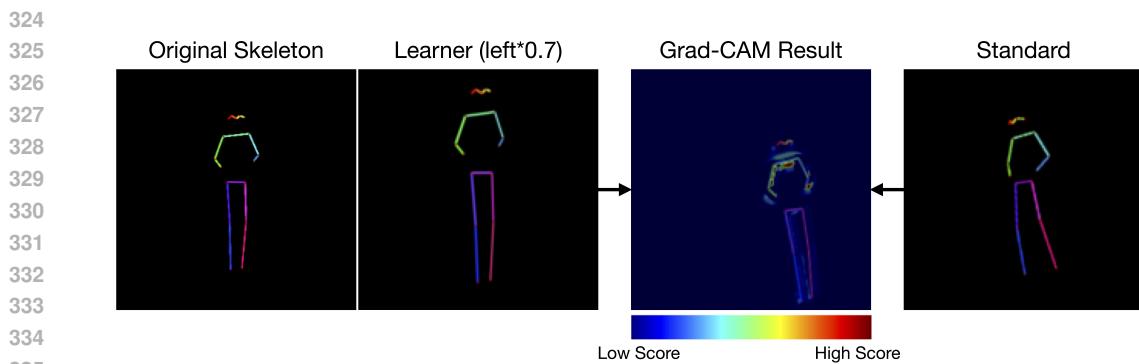


Figure 4: Example pairs of input poses and their corresponding Grad-CAM result

half and upper-half distortions reaches approximately 83%. This indicates that the baseline Siamese CNN tends to focus its activations on the upper body, even when the upper region is the one being distorted. As a result, the high accuracy observed for certain distortion directions (e.g., lower-half distortions) does not reliably reflect the model’s actual ability to identify disproportion; rather, it reflects a strong directional bias in its attention pattern.

In contrast, the STN-based model shows a substantially reduced gap—about 32%—between upper- and lower-half distortions, demonstrating that the STN mitigates this bias and produces more balanced activation patterns.

Compared with the left-right distortion results, the discrepancy between upper and lower distortions is particularly pronounced. In the axis-aligned baseline, Grad-CAM activations consistently concentrate in the upper part of the skeleton, regardless of whether the upper or lower half was distorted. This suggests that the model relies heavily on shoulder-level features when comparing skeletons, implying that these features change little under vertical stretching or compression. One likely explanation is that upper-body deformations occur in a direction perpendicular to the shoulder axis, and thus the shoulder structure remains relatively unaffected; as a consequence, the model cannot rely on this region to detect vertical proportional errors.

Moreover, more than two-thirds of our dataset consists of standing poses, in which the legs occupy a large visual area and align with the vertical direction. Because both upper- and lower-half distortions alter the relative leg length, the model becomes particularly sensitive to geometric changes in the lower body. Overlay visualizations further confirm that kneeling poses—whose leg segments do not exhibit a strong alignment with any single dominant axis—do not show the same activation concentration in the upper body. This reinforces the observation that the bias arises from the interaction between vertical deformations and the predominance of upright, vertically oriented poses in the dataset.

#### 4.5 EFFECT OF BENT SKELETONS

In the original dataset, all skeletons are strictly straight, which reflects typical usage scenarios but may limit the diversity of geometric patterns seen during training. To investigate whether the introduction of curvature can improve the model’s ability to extract robust skeletal features, we modify the dataset by removing the original swirl-based disproportion and adding four additional bending transformations. Each bending transformation is implemented using a swirl-like deformation with a different maximum rotation angle, namely 15°, 30°, 45°, or 60°. These bent skeletons account for approximately 18.2% of the total samples in the STN-based setting and 28.6% in the axis-aligned setting, and are treated as auxiliary training data rather than the primary focus.

For the STN-based model, incorporating bent skeletons keeps the overall accuracy at a similar level and even yields a slight improvement in some categories. Since these samples are processed by the STN along with all others, the model is encouraged to handle variations not only in size and position but also in local curvature. This additional variety helps the STN–Siamese CNN pipeline learn more

378 generalizable skeletal representations, and the model becomes better at extracting features that are  
 379 truly related to body proportions rather than tied to a single canonical pose.  
 380

381 In contrast, the axis-aligned baseline suffers a decrease in accuracy of about 5% after bent skeletons  
 382 are introduced. In this setting, bending is applied before the deterministic alignment of the central  
 383 axis. As a result, skeletons with strong curvature remain visually skewed even after alignment,  
 384 while the baseline model has no mechanism—unlike the STN-equipped model—to further adjust or  
 385 correct these residual distortions. These highly bent and misaligned examples act as outliers in the  
 386 training data and confuse the Siamese CNN, which in turn degrades its ability to judge proportional  
 387 correctness.

## 388 5 CONCLUSION

390 In this paper, we presented a proportionality-aware pose assessment framework for illustrated char-  
 391 acters and realized it through an STN–Siamese CNN architecture equipped with Grad-CAM visual-  
 392 ization. By constructing a balanced dataset that includes both proportion-preserving and proportion-  
 393 changing transformations, our method enables the STN to learn meaningful spatial alignment while  
 394 simultaneously guiding the Siamese CNN to capture structural differences between anatomically  
 395 correct and distorted skeletons. During inference, our system automatically matches a learner’s  
 396 skeleton to a correctly proportioned reference and generates intuitive heatmap-based feedback that  
 397 highlights correct regions while suppressing activations on distorted areas. Experimental results  
 398 demonstrate that our approach provides more balanced assessment compared with a deterministic  
 399 baseline, effectively reducing directional bias and maintaining performance even when additional  
 400 geometric variations. Given the lack of tools for structural guidance in anime-style illustration,  
 401 we believe that our proposed framework offers a promising direction for interpretable proportion  
 402 analysis and has significant potential for future development in educational and artistic applications.

## 404 6 WORKLOAD DISTRIBUTION

407 Table 2: Workload Distribution of Team Members

409 Name	410 Main Responsibilities
411 Si-Xian Hsu	Datasets construction, experiments design, model training, result analysis, paper and poster drafting
412 Yu-jyuan Lin	Survey, datasets construction, model training, model explo- ration, paper and poster drafting
413 Chian-ru Lin	Survey, model building, model training, paper and poster draft- ing

## 417 REFERENCES

- 419 Shuhong Chen and Matthias Zwicker. Transfer learning for pose estimation of illustrated characters.  
 420 In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 793–  
 421 802, 2022.
- 422 Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances  
 423 in neural information processing systems*, 28, 2015.
- 425 Yuhang Qiu, Jiping Wang, Zhe Jin, Honghui Chen, Mingliang Zhang, and Liquan Guo. Pose-  
 426 guided matching based on deep learning for assessing quality of action on rehabilitation training.  
 427 *Biomedical Signal Processing and Control*, 72:103323, 2022.
- 428 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,  
 429 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-  
 430 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,  
 431 2017.

432  
433

## A APPENDIX

434  
435

### A.1 LOW PERFORMANCES INTERPRETAION

436  
437

Although the overall accuracies of the experiments are not high, this outcome is not unexpected given that our application setting is substantially more challenging than the original rehabilitation setting, particularly in terms of dataset construction and the difficulty of the inference task. We propose several possible interpretations for these results.

440  
441  
442

The first factor is dataset complexity. Our dataset includes much higher diversity (different learner forms, and more than one standard example for each class) in order to reflect real art needs. This larger variance makes the model harder to train compared to the original rehabilitation dataset.

443  
444  
445

The second factor is inference task difficulty. Instead of comparing a learner’s pose to its unique standard, a learner must be compared to a random standard, because in realistic art practice there is no fixed standard pose. This makes it harder for the model to identify which parts have twisted from the original structure.

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485