PHI: Bridging Domain Shift in Long-Term Action Quality Assessment via Progressive Hierarchical Instruction

Kanglei Zhou, Hubert P. H. Shum, Senior Member, IEEE, Frederick W. B. Li, Xinxing Zhang, and Xiaohui Liang

Abstract-Long-term Action Quality Assessment (AQA) aims to evaluate the quantitative performance of actions in long videos. 2 However, existing methods face challenges due to domain shifts 3 between the pre-trained large-scale action recognition backbones 4 and the specific AQA task, hindering performance. This arises 5 since fine-tuning intensive backbones on small AQA datasets is 6 impractical. We address this by distinguishing domain shifts into task-level, regarding differences in task objectives, and 8 feature-level, regarding differences in important features. For 9 10 feature-level shifts, which are more detrimental, we propose Progressive Hierarchical Instruction (PHI) with two strategies. 11 First, Gap Minimization Flow (GMF) leverages flow matching 12 to progressively learn a fast flow path that reduces the domain 13 gap between initial and desired features across shallow to deep 14 layers. Additionally, a temporally-enhanced attention module 15 captures long-range dependencies essential for AOA. Second, 16 List-wise Contrastive Regularization (LCR) facilitates coarse-17 to-fine alignment by comprehensively comparing batch pairs to 18 learn fine-grained cues while mitigating shift. Integrating these, 19 PHI offers an effective solution. Experiments demonstrate that 20 PHI achieves state-of-the-art performance on three representative 21 long-term AQA datasets, proving its superiority in addressing the 22 domain shift issue for long-term AQA. 23

Index Terms—Action Quality Assessment, Long-Term Action 24 Quality Assessment, Domain Shift, Flow Matching 25

I. INTRODUCTION

26

Action Quality Assessment (AQA) [1], [2], [3], [4], [5] aims 27 to evaluate the quantitative performance of actions performed 28 in videos or image sequences. Unlike traditional action recog-29 nition, which focuses solely on identifying specific actions, 30 AQA provides a more detailed understanding of how well 31 those actions are executed [6]. This fine-grained evaluation 32 has broad applications across domains such as sports analysis 33

Manuscript received July 16, 2024; revised March 10, 2025. This work was supported by the National Natural Science Foundation of China under Project 62272019. (Corresponding author: Xiaohui Liang.)

Kanglei Zhou is with the State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing 100191, China, and also with the Department of Computer Science, Durham University, DH1 3LE Durham, U.K. (e-mail: zhoukanglei@buaa.edu.cn).

Hubert P. H. Shum and Frederick W. B. Li are with the Department of Computer Science, Durham University, DH1 3LE Durham, U.K. (e-mail: hubert.shum@durham.ac.uk; frederick.li@durham.ac.uk).

Xingxing Zhang is with the Department of Computer Science and Technology, Institute for AI, BNRist Center, Tsinghua-Bosch Joint ML Center, THBI Lab, Tsinghua University. (e-mail: xxzhang1993@gmail.com).

Xiaohui Liang is with the State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing 100191, China, and also with the Zhongguancun Laboratory, Beijing 100190, China (e-mail: liang_xiaohui@buaa.edu.cn).

¹Our code is included in the supplementary material for examination.



Fig. 1: Illustrations of our main idea: (a) The pre-trained I3D backbone emphasizes coarse features like guardrails (highlighted in vellow boxes), potentially unrelated to scoring for AQA, while it can accurately recognize cartwheeling in the action recognition domain. This discrepancy is primarily due to the pre-trained task's broader focus on coarse-level features, whereas fine-grained features essential for AQA may not be adequately exploited. (b) We identify two distinct domain shift types: task-level discrepancies and feature-level discrepancies. (c) Based on two hypotheses, our approach innovates a shallow-to-deep adaptation using Gap Minimization Flow (GMF), enabling a fast and controllable path to thoroughly minimize the domain gap. Additionally, we introduce a coarse-to-fine alignment mechanism using List-wise Contrastive Regularization (LCR) to enable the model to focus on fine-grained features, essential for AQA, while mitigating domain shift by refining coarse features from the broader pre-trained task.

[7], [8], [6], [9], medical rehabilitation [10], [11], and skill 34 assessment [12], [13], [14]. Long-term AQA [15], [16] extends AQA beyond individual snapshots or short clips to encompass extended durations. This broader evaluation is more challenging and practical, as it provides a comprehensive assessment 38 of actions in real-world scenarios.

One of the most significant challenges in long-term AQA is 40 the domain shift issue [17], where the pre-trained backbone is 41 suboptimal for AQA tasks. This challenge arises due to label 42 scarcity and the nature of long video sequences. Firstly, label 43 scarcity contributes to relatively small AQA datasets (e.g., one 44 large-scale long-term AQA dataset for rhythmic gymnastics 45 (balls) with approximately 250 samples). To address this, 46 most AQA methods [18], [8], [19] often leverage backbones 47 pre-trained on large-scale action recognition datasets (e.g., 48

35

36

37

Kinetics 400 [20] with over 300,000 samples). While this strat-1 egy enhances performance on small-scale AQA datasets, its 2 performance is restricted by the shift from action recognition 3 to AQA tasks (see Figure 1(a)). Secondly, the computational demands of processing these long sequences within long-term AQA, combined with the complexity of intensive backbones 6 [20], [21], make fine-tuning the backbone impractical with limited computation resources. As a result, existing long-term 8 AQA methods [16], [15] choose to fix the backbone and do not explicitly address the domain shift issue, thereby severely 10 limiting overall performance. Indeed, some methods [16], 11 [6] employing feature aggregation or representation layers, 12 such as Transformers [22], can implicitly mitigate domain 13 shift by leveraging score supervision to encourage the model 14 to capture high-level task-relevant patterns. However, these 15 methods remain ineffective for accurate assessment (refer to 16 results in Figures 7 and 8), as they lack explicit adaptation 17 mechanisms to align domain-specific feature distributions. 18

Long-term AQA aims to evaluate performance based on 19 important discriminative features, capturing subtle dynamics 20 over extended periods where domain shift significantly hinders 21 AQA performance. As illustrated in Figure 1(b), domain shift 22 arises from both differences in task objectives and variations in 23 data distribution. Accordingly, we categorize domain shift into 24 task-level discrepancies and feature-level discrepancies. Task-25 level discrepancies arise from transitioning from classification-26 based pre-training, where class boundaries are discrete, to 27 regression-based scoring for AQA, where scores vary con-28 tinuously. Feature-level discrepancies occur due to variations 29 in video capture conditions and application domains, leading 30 pre-trained models to focus on coarse, irrelevant patterns 31 rather than the fine-grained scoring cues necessary for precise 32 AQA. The previous work [23] addressed task-level discrep-33 ancies by formulating AQA as a coarse-to-fine classification 34 problem (see the alignment arrow "-" in the bottom of 35 Figure 1(b), but this resulted in precision loss. Our work 36 instead directly tackles feature-level discrepancies by refining 37 pre-trained coarse features to focus on fine-grained cues (see 38 " \rightarrow " in the top of Figure 1(b)). This refines pre-trained coarse 39 features to emphasize the fine-grained cues essential for ac-40 curate assessment without sacrificing precision. Experimental 41 results in Tables I to III demonstrate the superior performance 42 of our approach over addressing only task-level discrepancies, 43 showing the increased importance of minimizing feature-level 44 domain shifts for accurate long-term action assessment. 45

To this end, we introduce the Progressive Hierarchical 46 Instruction (PHI) framework (see Figure 1(c)) to tackle the 47 aforementioned challenges. Built on two major hypotheses 48 regarding shallow-to-deep adaptation and coarse-to-fine align-49 ment, we propose solutions to validate these hypotheses, 50 constituting the core components of PHI. These approaches 51 collectively reduce the domain gap while prioritizing fine-52 grained features essential for accurate assessment. 53

Hypothesis 1: A shallow-to-deep adaptation approach
through multiple-step control can thoroughly reduce the domain gap between a pre-trained backbone and the AQA task.
The rationale behind this lies in the complexity of refining
initial features in a single step, whereas the multi-step control

facilitates a more progressive and precise gap reduction across shallow to deep layers (see Figure 6).

To validate Hypothesis 1, we introduce Gap Minimization Flow (GMF) to progressively reduce the domain gap, which is motivated by the recent advances in flow models [24], [25]. However, these methods face challenges in constructing training pairs due to the inaccessibility of desired features. To address this, we first integrate temporally-enhanced attention to efficiently estimate desired features, capturing crucial longrange dependencies essential for long-term AQA. This enables GMF to directly and efficiently control the gap reduction, thereby enhancing the model's adaptability.

Hypothesis 2: A coarse-to-fine alignment approach that prioritizes learning representations focusing on fine-grained cues can effectively mitigate domain shift. The rationale behind this is that the pre-trained backbones on large-scale datasets often capture coarse features irrelevant to AQA scoring, whereas AQA depends on fine-grained cues (see Figure 1(a)).

To validate Hypothesis 2, we present List-wise Contrastive Regularization (LCR) to learn the fine-grained cues essential for AQA, which is motivated by the recent advances in contrastive regression [8], [9]. However, these methods are reliant on manual exemplar selection and known exemplar score during inference, restricting the application scope. To address this, LCR comprehensively compares all pairs of batch data, eliminating the need for manual intervention. In addition, this approach ensures a robust evaluation of action quality, especially in scenarios with limited labeled data.

Experimental results demonstrate the significant improvements achieved by PHI compared to state-of-the-art methods that do not specifically address domain shift issues. Notably, PHI achieves gains of 5.88%, 5.65%, and, 24.4% in correlation on three representative long-term AQA datasets, Rhythmic Gymnastics [15], Fis-V [26], and LOGO [3], respectively, compared to shift-unaware methods. Compared to the tasklevel solution [23], PHI showcases additional correlation gains and significant precision improvements, demonstrating the importance of addressing feature-level discrepancies for accurate long-term AQA. Our main contributions are as follows:

- We define domain shifts from both task and feature levels. Our work achieves additional performance gains in longterm AQA by addressing feature-level discrepancies.
- We propose a novel gap minimization flow module (GMF) to address the domain gap in a shallow-to-deep manner, facilitating efficient gap-reducing control with a fast and straight path.
- We design a novel contrastive regularization module (LCR) to mitigate the domain shift in a coarse-to-fine manner, enabling robust representation learning for performance improvement.

The remainder of this paper is organized as follows: Section II reviews the related work in AQA and flow matching methods; Section III briefly describes the preliminaries in rectified flow; Section IV details the core components of our proposed framework; Section V validates and analyzes the effectiveness of our proposed method; Section VI concludes the whole paper and outlooks the future work.

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

In this section, we provide a concise overview of AQA and
 flow matching, outlining their relevance to our work.

4 A. Action Quality Assessment (AQA)

1

AQA focuses on the quantitative performance of performed 5 actions in various application areas such as sports analy-6 sis [1], [7], [8], [6], [9], [23], medical rehabilitation [10], 7 [11], and skill assessment [12], [13], [14]. Earlier methods 8 depended heavily on hand-crafted features and heuristics, 9 revealing certain limitations. For example, Pirsiavash et al. [27] 10 leveraged pose features to train a linear SVR model, which 11 was constrained due to the poor pose estimation outcomes 12 [28]. By integrating deep learning, various models have shown 13 improved performance, such as CNNs [10], [6], RNNs [29], 14 [30], and Transformers [19], [23], [31]. 15

Existing AOA datasets [32] are relatively small, risking 16 over-fitting. To mitigate this, pre-trained backbones are com-17 monly employed. Parmar et al. [33] utilized C3D [34] to 18 improve AQA performance. Pan et al. [35] integrated a more 19 robust I3D backbone [20] to further optimize the performance. 20 Xu et al. [16] attempted to incorporate VST [21], aiming to 21 derive more powerful features. However, the computational 22 intensity of such 3D backbones presents a notable issue. As 23 every frame may contain essential AQA cues [26], videos 24 are often segmented into clips for separate processing, which 25 hinders a complete understanding of the action. To this end, 26 Zhou et al. [6] introduced a GCN method to eliminate semantic 27 ambiguities. Features extracted are typically aggregated either 28 by LSTM [29], or more popular average pooling before the 29 score regression. Instead of methods focusing on single-person 30 AQA datasets, Zhang et al. [3] proposed a group-aware diving 31 dataset. Furthermore, by leveraging the unique strengths of 32 different data types, such as video, audio, and skeleton data, 33 34 multi-modal AQA methods [36], [2], [37], [31] create a more holistic, accurate, and robust assessment system. In this work, 35 we primarily focus on RGB-based single-person AQA. 36

In the context of long-term AQA, the computational inten-37 sity of these backbones, combined with the long sequences, 38 poses significant challenges. Since fine-tuning the backbone 39 to minimize the domain shift between the pre-trained broader 40 task and the AQA task is difficult, existing long-term AQA 41 methods [16], [15] often choose to fix the backbone and over-42 look the domain shift problem, thereby limiting performance. 43 CoFInAl [23] addressed domain shifts by reformulating AQA 44 into coarse and fine classification tasks, potentially leading to 45 a loss in assessment precision. In contrast, our work identifies 46 and tackles the key challenge of domain shift through the lens 47 of feature-level discrepancies. Our method efficiently resolves 48 this issue without the need for fine-tuning the computational 49 backbone and preserving assessment precision. 50

51 B. Flow Matching

Flow Matching (FM) models represent a recent advancement in generative modeling. The term 'flow' refers to a mapping between samples of two distributions, leveraging neural Unlike traditional flow-based models, recent approaches propose training algorithms that require solving the ODE explicitly only during inference, overcoming challenges associated with backpropagating through ODEs during training [38], [25], [39], [24], [42]. FM offers a promising and littleexplored avenue for generative modeling, providing a more efficient and effective approach to learning complex distributions. Different from diffusion models [43], [44] that utilize Stochastic Differential Equations (SDEs) and are restricted to Gaussian base distributions [45], FM offers greater flexibility by allowing the choice of base distribution and training with ODEs instead of SDEs, leading to smoother trajectories and improved performance [46].

Our work is inspired by the rectified flow [24] that generates a fast and direct flow path but relies on known target samples, which poses significant challenges in addressing domain shift. To overcome this limitation, we introduce a novel approach that does not depend on any known target distributions. We emphasize that our proposed PHI method is not merely an adjustment to rectified flow but a fundamental extension that introduces a novel hierarchical adaptation framework. By leveraging shallow-to-deep adaptation and coarse-to-fine alignment, PHI achieves self-supervised domain adaptation without explicit target distribution samples, making it highly generalizable to a wide range of real-world applications beyond AQA. Notably, PHI is the first attempt at integrating the concept of flow matching in the realm of AQA.

III. PRELIMINARY: RECTIFIED FLOW

Rectified Flow (RF) [24] offers a straightforward solution for finding a transport map between two observed distributions. It involves learning an Ordinary Differential Equation (ODE), also known as the flow model, aiming to traverse straight paths as much as possible. Given empirical observations $\boldsymbol{x}_0 \sim \pi_0, \boldsymbol{x}_1 \sim \pi_1$, RF finds the transport map implicitly by solving the following ODE problem:

$$d\boldsymbol{z}_t = v(\boldsymbol{z}_t, t)dt, \quad t \in [0, 1], \tag{1}$$

where $v: \mathbb{R}^d \to \mathbb{R}^d$ represents the drift force that converts z_0 from distribution π_0 to z_1 following distribution π_1 .

RF operates by aligning the ODE with the linear inter-98 polation of points from π_0 and π_1 . Given x_0 and x_1 , the 99 linear interpolation of x_0 and x_1 is $x_t = tx_1 + (1-t)x_0$. 100 Observe \boldsymbol{x}_t follows a trivial ODE that already transfers π_0 101 to π_1 , i.e., $\mathrm{d} \boldsymbol{x}_t = (\boldsymbol{x}_1 - \boldsymbol{x}_0) \mathrm{d} t$, where \boldsymbol{x}_t moves following 102 the line direction $(x_1 - x_0)$ with a constant speed. However, 103 this ODE does not solve the problem: it can't be simulated 104 causally, because the update x_t depends on the final state x_1 , 105 which is not supposed to be known at time t > 1. 106

RF casualizes the interpolation process by projecting it to ¹⁰⁷ the space of causally simulatable ODEs in Equation (1). The ¹⁰⁸

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96



Fig. 2: Framework of PHI: Our PHI framework addresses the domain shift issue through two crucial processes. Firstly, Gap Minimization Flow (GMF) progressively transforms the initial feature into the desired one, minimizing the domain gap. Secondly, List-wise Contrastive Regularization (LCR) guides the model towards subtle variations in actions, facilitating the transition from coarse to fine-grained features crucial for AQA. Finally, the refined feature is used to predict the quality score through an MLP.

¹ drift force v is set to drive the flow to follow the direction ² $(x_1 - x_0)$ as much as possible. A natural way to the ℓ_2 ³ projection on the velocity field is finding v by solving a simple ⁴ least squares regression problem:

$$\min_{v} \int_{0}^{1} \mathbb{E}\left[\|(\boldsymbol{x}_{1} - \boldsymbol{x}_{0}) - v(\boldsymbol{x}_{t}, t)\|^{2} \right] \mathrm{d}t.$$
 (2)

⁵ By fitting v with the direction $(x_1 - x_0)$, RF casualizes the ⁶ paths of linear interpolation x_t , yielding an ODE flow that can ⁷ be simulated without seeing the future. We can parameterize ⁸ v with a neural network and solve Equation (2) with any ⁹ stochastic optimizer.

While RF aims to find a transport map between two ob-10 served distributions by projecting the interpolation path onto 11 a space of causally simulatable ODEs, our proposed Gap Min-12 imization Flow module extends this concept to progressively 13 minimize the domain gap between the pre-trained backbone 14 and the target AQA task. Unlike the ODE-based formulation 15 in RF, our approach leverages temporally-enhanced attention 16 to efficiently estimate the desired features, allowing for a more 17 targeted and effective domain gap reduction. 18

IV. METHODOLOGY: PHI

This section first introduces the PHI framework, followed by a detailed explanation of its core components.

19

a) Problem Definition: Long-term AQA involves assess-22 ing the quality or proficiency of actions and activities through 23 extended video recordings, considering temporal factors such 24 as consistency, progression, and the detailed evolution of 25 the full performance over time, rather than just snapshots. 26 This presents significant challenges in addressing the domain 27 shift issue, modeling long-range temporal dependencies, and 28 capturing fine-grained dynamics. 29

b) Framework Overview: Our PHI method (see Figure 2) addresses the aforementioned challenges through a holistic two-stage methodology involving shallow-to-deep adaptation and coarse-to-fine alignment. Given an action video $\mathbf{X}_i \in \mathbb{R}^{T \times W \times H \times 3}$, representing T frames of resolution $W \times H$ and 3 color channels, the initial feature $\mathbf{H}_i^0 \in \mathbb{R}^{M \times D}$ is extracted using a pre-trained backbone, where M denotes the number 36 of clips. In the shallow-to-deep adaptation stage (see Sec-37 tion IV-A), Gap Minimization Flow (GMF) with temporally-38 enhanced attention transforms the initial feature \mathbf{H}_{i}^{0} into the 39 desired feature representation $\mathbf{H}_{i}^{1} \in \mathbb{R}^{M \times D}$. The loss function 40 \mathcal{L}_{M} (see Equation (11)) facilitates efficient adjustment of the 41 initial feature to achieve the desired representation that better 42 aligns with AQA. In the coarse-to-fine alignment stage (see 43 Section IV-B), List-wise Contrastive Regularization (LCR) 44 regularizes the feature space. The regularization loss $\mathcal{L}_{\rm R}$ (see 45 Equation (15)) encourages the adaptation process to prioritize 46 fine-grained features essential for AOA. Ultimately, the desired 47 feature representation \mathbf{H}_{i}^{1} is fed into an MLP head network to 48 predict the quality score \hat{s}_i , which is supervised by the MSE 49 loss $\mathcal{L}_{\rm S} = \frac{1}{2} \sum_{i} (s_i - \hat{s}_i)^2$. Overall, the total loss is: 50

$$\mathcal{L} = \mathcal{L}_{\rm S} + \lambda_{\rm M} \mathcal{L}_{\rm M} + \lambda_{\rm R} \mathcal{L}_{\rm R}, \qquad (3)$$

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

where λ_{M} and λ_{R} are loss weights. During testing, the initial feature only undergoes the flow path to regress the score.

A. GMF: Shallow-to-Deep Adaptation

1)

Hypothesis 1: A shallow-to-deep adaptation approach through multiple-step control can thoroughly reduce the domain gap between a pre-trained backbone and the AQA task.

a) Justification: The challenge of mitigating the domain shift arises from the complexity of transferring complex data between initial and desired features in a single step, often compromising reliability and precision. A multi-step control mechanism involves breaking down the distribution gap into multiple sub-parts and sequentially matching these distributions. This approach enables gradual and accurate feature transformation, facilitating a thorough reduction of the domain gap across shallow to deep layers.

b) Implementation: The core idea of GMF is to utilize the concept of RF [24] to gradually reduce the domain gap, which relies on coupling pairs to train the flow path. However, in our context, we lack the corresponding desired feature for



Fig. 3: Illustration of Temporally-Enhanced Transformer Encoder (TETE): Highlighting the attention of Temporally-Enhanced Self-Attention (TESA), we employ the low-rank matrix to reduce the complexity from $\mathcal{O}(M^2D)$ to $\mathcal{O}(Md_tD)$, ensuring efficient modeling of long-term dependencies.

the initial feature, posing a significant challenge in minimizing the domain gap. To address this, we first propose to estimate the desired feature and then progressively minimize the gap. Fortunately, the semantics of the desired feature are known, and represented by the score. Thus, leveraging score loss $\mathcal{L}_{\rm S}$ allows for the coarse estimation of the desired feature. Our shallow-to-deep adaptation consists of two steps: desired feature estimation and domain gap minimization.

2) Desired Feature Estimation: In long-term AQA, effectively modeling long-range dependencies in sequences is 10 crucial. The previous work [23] employs a simple transforma-11 tion matrix to aggregate temporal dependencies, which, while 12 effective, cannot capture long-range temporal relationships. 13 While self-attention can model the long-range dependencies, it 14 entails substantial computational overhead. To address this, we 15 propose Temporally-Enhanced Transformer Encoder (TETE, 16 see Figure 3) using low-rank decomposition to optimize com-17 putation within attention. TETE can efficiently model long-18 range dependencies in long-term AOA by reducing com-19 putational complexity in attention mechanisms and enabling 20 accurate estimation of desired features. 21

22 a) Vanilla Self-Attention: Considering the initial feature 23 $\mathbf{H}_{i}^{0} \in \mathbb{R}^{M \times D}$ of *i*-th action video, the vanilla self-attention 24 [22] can be represented as:

Attention(
$$\mathbf{Q}, \mathbf{K}, \mathbf{V}$$
) = Softmax $\left(\frac{\mathbf{QK}^{\top}}{\sqrt{d_k}}\right) \mathbf{V}$, (4)

where $\mathbf{Q} \in \mathbb{R}^{M \times d_k}$, $\mathbf{K} \in \mathbb{R}^{M \times d_k}$, $\mathbf{V} \in \mathbb{R}^{M \times D}$ are the linear embeddings of the input \mathbf{H}_i^0 . The computation complexity of Equation (4) is $\mathcal{O}(M^2D)$, where the clip number M is typically large for long-term AQA. Thus, computation increases sharply with longer sequences in real-world scenarios.

b) Temporally-Enhanced Self-Attention (TESA): Our 30 proposed TESA module addresses the computational chal-31 lenges of long sequences by learning a low-rank matrix 32 $\in \mathbb{R}^{d_t \times d_k}$, where d_t is significantly smaller than the \mathbf{T} 33 feature dimension D, ensuring computational efficiency while 34 preserving temporal dependencies. First, TESA applies T only 35 to the query \mathbf{Q} and the key \mathbf{K} , as they directly determine 36 the attention weights. This enhances the ability of the model 37 to capture long-term dependencies while avoiding redundant 38

transformations on the value **V**, which does not influence attention weight computation. The transformed queries and keys are computed as: 41

$$\mathbf{T}_q = \operatorname{Softmax} \left(\mathbf{Q} \mathbf{T}^\top \right) \in \mathbb{R}^{M \times d_t}, \tag{5}$$

$$\mathbf{\Gamma}_k = \operatorname{Softmax}(\mathbf{K}\mathbf{T}^\top) \in \mathbb{R}^{M \times d_t}.$$
(6)

Next, we use T_k to query the value V, yielding:

$$\mathbf{\Gamma}_{v} = \operatorname{Softmax}(\mathbf{T}_{k}\mathbf{V}) \in \mathbb{R}^{d_{t} \times D}.$$
(7)

Finally, our attention mechanism can be represented as:

Attention
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{T}) = \mathbf{T}_q \mathbf{T}_v \in \mathbb{R}^{M \times D}$$
. (8)

This results in a decent complexity reduction to $\mathcal{O}(Md_tD)$, effectively improving the training efficiency and mitigating overfitting, thus leading to performance improvement (refer to results in Table V).

Finally, the desired feature can be estimated as:

$$\mathbf{H}_{i}^{1} = \mathcal{T}\left(\mathbf{H}_{i}^{0}\right),\tag{9}$$

where $\mathcal{T}(\cdot)$ represents the TETE module.

3) Domain Gap Minimization: Inspired by the concept of RF [24], our work aims to learn a fast and controllable path to efficiently mitigate domain shift by progressively reducing the domain gap between initial and desired features.

Suppose the initial feature \mathbf{H}_{i}^{0} and the desired feature \mathbf{H}_{i}^{1} 54 follow two different empirical distributions, and we sample 55 P intermediate steps. The target representation of the j-56 th step can be defined as the linear interpolation $\mathbf{H}_{i}^{j/F}$ = 57 $\frac{j}{D}\mathbf{H}_{i}^{0} + (1 - \frac{j}{D})\mathbf{H}_{i}^{1}$. During training, our objective is to predict 58 the next step using the previous representation and the step 59 size through a neural network $\phi(\cdot, \cdot)$, indicating the domain 60 gap of the corresponding step. In practice, our experiments 61 demonstrate that a simple MLP is sufficient for this task. The 62 gap at the *j*-th step can be represented as: 63

$$\boldsymbol{g}_j = \phi\left(\hat{\mathbf{H}}_i^{(j-1)/P}, \frac{1}{P}\right), \quad j \ge 1,$$
(10)

where $\hat{\mathbf{H}}_{i}^{(j-1)/P}$ denotes the predicted feature of the previous step, calculated as $g_{j-1} + \hat{\mathbf{H}}_{i}^{(j-2)/P}$ if $j \ge 2$, otherwise \mathbf{H}_{i}^{0} . The network ϕ generates a driving force that guides the flow in the direction of the overall flow, and we optimize ϕ by: 67

$$\mathcal{L}_{\mathrm{M}} = \frac{1}{B} \sum_{i=0}^{B-1} \left(\mathcal{L}_{\mathrm{M-global}} + \mathcal{L}_{\mathrm{M-local}} \right), \qquad (11)$$

$$\mathcal{L}_{\mathrm{M-global}} = \left\| \left(\mathbf{H}_{i}^{1} - \mathbf{H}_{i}^{0} \right) - \sum_{j=1}^{P} g_{j} \right\| , \qquad (12)$$

$$\mathcal{L}_{\mathrm{M-local}} = \frac{1}{P} \sum_{j=1}^{P} \left\| \mathbf{H}_{i}^{j/P} - \hat{\mathbf{H}}_{i}^{j/P} \right\|^{2}, \tag{13}$$

where *B* denotes the batch size, and $(\mathbf{H}_{i}^{1} - \mathbf{H}_{i}^{0})$ represents the overall flow direction. The first term $\mathcal{L}_{M-global}$ represents the global flow constraint, while the second term $\mathcal{L}_{M-local}$ denotes the local flow constraint. These two terms are collectively effective in addressing different aspects of the flow constraints.

42

43

45 46 47

48

49

50

51

52

53

4) Benefits of GMF: GMF introduces a fundamentally new perspective on feature transformation in flow-based adaptation, 2 offering several key advantages over existing methods. Unlike з traditional flow methods [25], [24], which require explicit access to the target distribution, GMF eliminates this dependency by leveraging an adaptive hierarchical transformation 6 strategy. This enables GMF to handle domain-shifted tasks where the target distribution is either unknown or difficult 8 to obtain, significantly extending the applicability of flowbased learning. Instead of directly aligning source and target 10 distributions, GMF sequentially matches intermediate rep-11 resentations through a multi-stage adaptation process. This 12 hierarchical alignment ensures a smooth and stable transi-13 tion while preventing abrupt transformations. Unlike stacked 14 or recurrent architectures, GMF dynamically regulates the 15 adaptation trajectory, reducing computational overhead and 16 improving convergence efficiency. Additionally, its implicit 17 knowledge distillation mechanism allows for lightweight infer-18 ence without compromising performance. These innovations 19 make GMF a scalable and efficient solution for domain-20 adaptive learning in video-based tasks, surpassing conventional 21 flow-based approaches. 22

B. LCR: Coarse-to-Fine Alignment 23

1) Design Idea: 24

56

Hypothesis 2: A coarse-to-fine alignment approach that 25 prioritizes learning representations focusing on fine-grained 26 cues can effectively mitigate the domain shift issue in long-27 term AQA. 28

a) Justification: As shown in Figure 1(a), the pre-trained 29 broader task often emphasizes coarse-level features that may 30 not directly align with the requirements of AQA. Conversely, 31 AQA tasks rely heavily on fine-grained cues [47], [48] for 32 accurate assessment. Therefore, a method that prioritizes learn-33 ing these fine-grained representations is likely to be more 34 effective in addressing the domain shift issue. 35

b) Implementation: Notably, the interpretation of coarse-36 to-fine alignment differs from CoFInAl [23]. In CoFInAl, 37 coarse-to-fine alignment refers to structuring AQA as a two-38 stage classification task, which helps address task-level differ-39 ences but may introduce precision loss. In contrast, PHI adopts 40 a feature refinement approach that progressively enhances pre-41 trained coarse features to focus on fine-grained scoring cues, 42 ensuring improved AQA performance. In response, our novel 43 LCR module (see Figure 4) employs a coarse-to-fine alignment 44 strategy to further address the domain shift, enhancing the 45 model's ability to capture fine-grained features essential for 46 AQA. By comparing the differences between actions, com-47 parative regression aids in identifying subtle variations or 48 abnormalities that may not be apparent when assessing actions 49 in isolation. Given a batch of data, LCR involves computing 50 a distance matrix where each row encodes the relationships 51 of an action with all other actions in the batch. By aligning 52 the distribution of each row with its ground truth quality score 53 distribution, the model can effectively learn subtle variations. 54 2) Action Distance Computation: For each action in the 55 batch, we need to calculate its pairwise distance with all other 69

70

71

72

73

74

75

76

77

78

79

80

81

82



Fig. 4: Illustration of List-wise Contrastive Regularization (LCR): LCR aligns the distance distribution in the feature space with that of the quality score space in a list-wise manner. This ensures comprehensive comparison and alignment across the entire batch of data, leading to robust performance.

actions. However, directly measuring the distance between 57 two actions is inappropriate because their clips may not be 58 temporally or semantically aligned. For instance, clips with 59 similar semantics might appear at different indices in each 60 action sequence. Hence, we need to establish correspondence 61 between clips before computing the action-level distance. (1) 62 To align clips in two actions, For each clip in the first action, 63 we compute its distance to all clips in the second action and 64 select the closest one as its paired clip. Specifically, the pairing 65 is determined based on the minimum ℓ_2 distance. (2) Then, 66 we sum all the clip pair distances to obtain the final distance 67 between the two actions. This process can be represented as: 68

$$D_{ij} = \sum_{m=0}^{M-1} \left(\min_{n \in [0, M-1]} (d_{ij}^{mn}) \right), \ d_{ij}^{mn} = \| \boldsymbol{h}_i^m - \boldsymbol{h}_j^n \|^2, \ (14)$$

where d_{ij}^{mn} denotes the distance between the *m*-th clip of the *i*-th action video and the n-th clip of the *j*-th action video using the ℓ_2 distance.

It is interesting to note the implicit temporal order relationship between the paired clip distance calculation. In practice, when considering two clips $h_i^{m_1}$ and $h_i^{m_2}$ $(m_1 \leq m_2)$ from the same video, where m_1 and m_2 represent different time steps, and their paired clips $m{h}_j^{n_1}$ and $m{h}_j^{n_2}$, it is observed that upon convergence, n_1 tends to be less than or equal to n_2 . This observation aligns with the inherent temporal ordering within action sequences, where paired clips that are closer in time tend to have smaller indices than those further apart. Therefore, when we added a loss item to constrain $n_1 \leq n_2$, the performance did not change.

3) Distance Matrix Alignment: Aligning the distribution of 83 each row in the distance matrix with the ground truth quality 84 score distribution involves two main steps. (1) We calculate 85 the ground truth quality score distance using $\mathbf{S} = |s - s^{\top}|$, 86 where $s \in \mathbb{R}^{B \times 1}$ denotes the quality score of the batch data 87 (B is the batch size). (2) Optimizing the alignment process 88 can be challenging, as traditional metrics like MSE may not 89 accurately capture the underlying structure of the data (see 90 results in Table V). To address this issue, we propose incor-91 porating Kullback-Leibler (KL) divergence [49] to enhance the 92 alignment process and better capture the complex relationships 93 J

within the data distributions. This process can be optimized by:

$$\mathcal{L}_{\mathrm{R}} = \sum_{i=0}^{B-1} \left(\mathrm{KL}(\boldsymbol{S}_i \| \boldsymbol{D}_i) + \mathrm{KL}(\boldsymbol{D}_i \| \boldsymbol{S}_i) \right) \right).$$
(15)

Here, the adoption of a symmetrized and smoothed version 2 of the KL divergence enhances the robustness and stability 3 of the alignment process, leading to accurate assessment. The 4 estimated target features generated by TETE (supervised by 5 \mathcal{L}_{S}) may not fully eliminate the domain shift. Therefore, GMF 6 aligns initial features with the estimated ones and then realigns these estimated features with the ideal ones essential for 8 AQA. LCR (with \mathcal{L}_R) regularizes this re-alignment, enhancing 9 overall performance. 10

4) Benefits of LCR: LCR introduces a novel hierarchi-11 cal alignment strategy that significantly enhances AQA per-12 formance by capturing subtle relationships between actions 13 within a batch. Unlike traditional pair-wise contrastive meth-14 ods [50], [8], [47], LCR employs batch-wise contrastive learn-15 ing to explicitly consider all pair-wise differences, enabling 16 a more effective capture of fine-grained variations. The loss 17 function \mathcal{L}_{R} leverages KL divergence for efficient alignment 18 between distributions D_i and S_i , ensuring faster convergence 19 and better generalization compared to direct alignment meth-20 ods like MSE. Additionally, low-rank regularization enforces 21 temporal coherence and robustness, preventing overfitting to 22 spurious correlations. The batch-wise learning and low-rank 23 constraints collectively establish LCR as a theoretically unique 24 and generalizable framework for domain adaptation, applicable 25 not only to AQA but also to tasks like rehabilitation analysis, 26 sports motion scoring, and action recognition. By addressing 27 the limitations of prior methods, LCR represents a significant 28 advancement in feature alignment for temporal and domain-29 shifted tasks. 30

31

V. EXPERIMENTS

In this section, we first describe the experimental setup, and then present and analyze the experimental results.

34 A. Experimental Setups

a) Datasets: In our study, we assess all models using 35 three extensive long-term AQA datasets. The first dataset, 36 named the Rhythmic Gymnastics (RG) dataset [15], com-37 prises a collection of 1,000 videos showcasing various rhyth-38 mic gymnastics actions performed with different apparatuses, 39 including Ball, Clubs, Hoop, and Ribbon. Each video has an 40 approximate duration of 1.6 minutes, and the frame rate is 41 set at 25 frames per second. The dataset is split into training 42 and evaluation sets, with 200 videos allocated for training 43 and 50 for evaluation in each action category. The second 44 dataset, referred to as the Figure Skating Video (Fis-V) 45 dataset [27], [26], contains 500 videos capturing ladies' singles 46 short programs in figure skating. Each video has a duration 47 of approximately 2.9 minutes, with a frame rate set at 25 48 frames per second. Following the official split, the dataset is 49 divided into 400 training videos and 100 testing videos. Each 50 video in this dataset comes with annotations for two scores: 51

Total Element Score (TES) and Total Program Component Score (PCS). To align with the previous method [29], we develop separate models for different score/action types. The third dataset, i.e., the **LOng-form GrOup** (**LOGO**) dataset [3], consists of 150 samples for training and 50 for testing. It captures videos showcasing synchronized swimming group actions, where each video sequence is approximately 3 and a half minutes in length. Currently, LOGO has the longest video duration among all AQA datasets, making it a challenging benchmark for long-term AQA tasks.

b) Evaluation Metrics: We use two evaluation metrics to validate the performance of all the AQA methods.

Consistent with previous long-term AQA methods [16], [15], we utilize Spearman's Rank Correlation Coefficient (SRCC) as the evaluation metric, denoted as ρ . The SRCC is defined as the Pearson correlation coefficient between their ranks, $r(s_i)$ and $r(\hat{s}_i)$, to predicted and ground-truth scores, which can be formulated as follows:

$$\rho = \frac{\sum_{i=1}^{N} \left(r\left(s_{i}\right) - \bar{r} \right) \left(r\left(\hat{s}_{i}\right) - \bar{r} \right)}{\sqrt{\sum_{i=1}^{N} \left(r\left(s_{i}\right) - \bar{r} \right)^{2}} \sqrt{\sqrt{\sum_{i=1}^{N} \left(r\left(\hat{s}_{i}\right) - \bar{r} \right)^{2}}}, \quad (16)$$

where \bar{r} is the average rank. A higher SRCC indicates a stronger rank correlation between predicted and ground-truth scores. Following the previous work [35], we compute the average SRCC across different action types for the RG dataset and score types for the Fis-V dataset by aggregating individual SRCCs using Fisher's z-value.

Compared with the previous work [23], we have added a new stricter metric, the relative ℓ_2 distance $(R-\ell_2)$ [8], [6]. The purpose of introducing $R-\ell_2$ is to measure the relative error of AQA models more precisely without being affected by the score scale. Given the highest and lowest scores for an action s_{max} and s_{min} , the relative ℓ_2 distance $R-\ell_2$ is defined as:

$$\mathbf{R} \cdot \ell_2 = \frac{1}{N} \sum_{n}^{N} \left(\frac{|s_n - \hat{s}_n|}{s_{\max} - s_{\min}} \right)^2 \times 100, \qquad (17)$$

where s_n and \hat{s}_n represent the ground-truth score and prediction for the *n*-th sample, respectively. Fisher's z-value is used to measure the average performance across actions.

c) Implementation Details: We implemented PHI using 85 PyTorch on an RTX 3090 GPU. We employ VST pre-trained 86 on Kinetics 600 [16] and I3D on Kinetics 400 [20] as the 87 backbone to conduct experiments, respectively. The feature 88 dimensions D, d_k, d_t are set to 1024, 128, and 32, respectively. 89 Following the previous work [16], [3], we initially partition 90 the video into non-overlapping 32-frame segments. During 91 training, we randomly determine the start segment, specifically 92 M = 68 for RG, M = 124 for Fis-V, and 48 for LOGO, 93 respectively. During testing, all segments are utilized. We 94 optimize all models using SGD with a momentum of 0.9. 95 The batch size B is 32, and the learning rate starts at 0.01, 96 gradually decreasing to 0.0001 through a cosine annealing 97 strategy. For convergence, all the models are trained for 200 98 epochs. The loss weights $\lambda_{\rm M}, \lambda_{\rm R}$ are set to 0.5 and 0.01, 99 respectively. To further optimize the networks, we apply a 100 dropout of 0.3 and a weight decay of 0.01. 101

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

82

83

TABLE I

Results of SRCC (\uparrow) and R- ℓ_2 (\downarrow) on the RG dataset. The best results are highlighted in **BOLD**, while the second-best results are UNDERLINED. THE SYMBOL "*" INDICATES OUR REIMPLEMENTATION BASED ON THE OFFICIAL CODE. THE AVERAGE SRCC IS COMPUTED USING THE FISHER-Z VALUE. "-" DENOTES THAT THE METHOD DOES NOT REPORT THIS METRIC. "+" INDICATES THE USE OF ADDITIONAL FEATURES/MODALITIES.

	D 1 V 1	hilling Dealthana			Ba	11	Ch	ubs	Но	op	Rib	bon	Ave	rage
Method	Publisher	Backbone	Aware	Modality	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$
C3D+SVR [26]	TPAMI'17	C3D	×	RGB	0.357	-	0.551	-	0.495	-	0.516	-	0.483	-
MS-LSTM [29]	TCSVT'19	I3D	×	RGB	0.515	-	0.621	-	0.540	-	0.522	-	0.551	-
ACTION-NET [15]	ACM MM'20	I3D+	×	RGB	0.528	-	0.652	-	0.708	-	0.578	-	0.623	-
GDLT [16]	CVPR'22	I3D	X	RGB	0.526	2.943	0.710	2.557	0.729	8.149	0.704	3.485	0.674	4.284
HGCN* [6]	TCSVT'23	I3D	×	RGB	0.534	6.748	0.609	16.142	0.706	9.270	0.621	9.934	0.621	10.524
VATP-Net [51]	TCSVT'24	I3D	×	RGB+	0.580	-	0.720	-	0.739	-	0.724	-	0.709	-
CoFInAl [23]	IJCAI'24	I3D	1	RGB	0.625	2.647	0.719	3.093	0.734	3.892	0.757	2.607	0.712	3.060
PHI (Ours)	-	I3D	1	RGB	0.598	3.471	0.732	3.139	0.731	5.376	0.754	5.674	0.708	4.415
MS-LSTM [29]	TCSVT'19	VST	×	RGB	0.621		0.661	-	0.670		0.695	-	0.663	-
ACTION-NET [15]	ACM MM'20	VST+	×	RGB	0.684	-	0.737	-	0.733	-	0.754	-	0.728	-
GDLT [16]	CVPR'22	VST	X	RGB	0.746	2.833	0.802	2.179	0.765	2.012	0.741	2.579	0.765	2.401
HGCN* [6]	TCSVT'23	VST	×	RGB	0.711	3.030	0.789	3.444	0.728	5.312	0.703	5.576	0.735	4.341
PAMFN [31]	TIP'24	VST	Х	RGB	0.636	-	0.720	-	0.769	-	0.708	-	0.711	-
VATP-Net [51]	TCSVT'24	VST	×	RGB+	0.800	-	0.810	-	0.780	-	0.769	-	0.800	-
CoFInAl [23]	IJCAI'24	VST	1	RGB	0.809	1.356	0.806	2.453	0.804	9.918	0.810	2.383	0.807	4.028
PHI (Ours)	-	VST	1	RGB	0.818	<u>2.187</u>	<u>0.803</u>	2.149	0.812	<u>2.119</u>	0.805	2.744	0.810	2.300

TABLE II

Results of SRCC (\uparrow) and R- ℓ_2 (\downarrow) on the Fis-V dataset. The best results are highlighted in **Bold**, while the second-best results are UNDERLINED. THE SYMBOL "*" INDICATES OUR REIMPLEMENTATION BASED ON THE OFFICIAL CODE. THE AVERAGE SRCC IS COMPUTED USING THE FISHER-Z VALUE. "-" DENOTES THAT THE METHOD DOES NOT REPORT THIS METRIC. "+" INDICATES THE USE OF ADDITIONAL FEATURES/MODALITIES.

			Shift Mark		TES		PCS		Average	
Method	Publisher	Backbone	Aware	Modality	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$
C3D+SVR [26]	TPAMI'17	C3D	×	RGB	0.400	-	0.590	-	0.501	-
MS-LSTM [29]	TCSVT'19	C3D	×	RGB	0.650	-	0.780	-	0.721	-
GDLT [16]	CVPR'22	I3D	×	RGB	0.260	5.582	0.395	5.039	0.329	5.311
HGCN* [6]	TCSVT'23	I3D	×	RGB	0.311	4.317	0.407	4.608	0.360	4.463
CoFInAl [23]	IJCAI'24	I3D	1	RGB	0.589	3.470	0.788	2.843	0.702	3.157
PHI (Ours)	-	I3D	1	RGB	0.659	2.572	0.798	<u>3.073</u>	0.736	2.823
MS-LSTM [29]	TCSVT'19	VST	×	RGB	0.660	-	0.809	-	0.744	-
ACTION-NET [15]	ACM MM'20	VST+	×	RGB	0.694	-	0.809	-	0.757	-
GDLT [16]	CVPR'22	VST	X	RGB	0.685	3.717	0.820	2.072	0.761	2.895
HGCN* [6]	TCSVT'23	VST	×	RGB	0.246	12.628	0.221	20.531	0.234	16.580
MLP-Mixer [37]	AAAI'23	VST	X	RGB	0.680	-	0.820	-	0.750	-
SGN [36]	TMM'24	VST	×	RGB	0.700	-	0.830	-	0.765	-
PAMFN [31]	TIP'24	VST	X	RGB	0.665	-	0.823	-	0.755	-
VATP-Net [31]	TCSVT'24	VST	×	RGB+	0.702	-	0.863	-	0.796	-
CoFInAl [23]	IJCAI'24	VST	1	RGB	0.716	<u>2.875</u>	0.843	1.752	0.788	2.314
PHI (Ours)	-	VST	1	RGB	0.726	2.543	0.867	1.656	0.804	2.178

B. Results and Analysis

We begin with a thorough comparison to state-of-the-art methods, followed by an ablation study that includes parameter 3 sensitivity analysis, and conclude with visualizations.

1) Comparison with the State-of-the-Art: In our compar-5 ative analysis, we benchmarked several state-of-the-art meth-6 ods, consisting of C3D+SVR [26], MS-LSTM [29], ACTION-NET [15], GDLT [16], HGCN [6], and CoFInAl [23]. The comprehensive results are reported in Tables I to III, and the 9 computational performance is listed in Table IV. 10

a) Comparison with Different Backbones: We employed 11 backbone various architectures, namely C3D [34], I3D [20], 12 ResNet [52], and VST [21], to discern their impact on AQA 13 performance. Among them, I3D consistently outperformed 14 C3D across all categories on the RG dataset, highlighting 15 its proficiency in capturing spatial-temporal dynamics for 16 AQA. Particularly noteworthy was the superior performance 17 of the VST backbone, achieving the highest average SRCC 18

on all the three datasets, as can be seen in Tables I to III. 19 In the RG dataset, ACTION-NET, when coupled with the 20 VST backbone, displayed a remarkable correlation gain of 21 over 16.85% compared to its I3D-based counterpart. This 22 highlights the VST backbone's capability to capture detailed 23 temporal dynamics crucial for AQA. Our proposed method, 24 PHI, leveraging both I3D and VST backbones, showcased 25 outstanding results across all categories, with the VST variant 26 emerging as the top performer. These findings validate the 27 significance of advanced 3D convolutional architectures, such 28 as the I3D and VST backbones, in capturing subtle action 29 details for improved AQA performance, further enhanced by 30 the incorporation of our PHI method. 31

b) Comparison with Shift-Unaware Methods: Traditional 32 shift-unaware AQA methods [16], [6], [29], [15] employ 33 representation layers that may implicitly mitigate the domain shift issue. However, the persistence of domain shift often leads to suboptimal performance in Tables I to III. PHI 36 with VST consistently outperforms others across all actions

34

35

TABLE III Results of SRCC (\uparrow) and R- ℓ_2 (\downarrow) on the LOGO dataset. The best RESULTS ARE HIGHLIGHTED IN BOLD, WHILE THE SECOND-BEST RESULTS ARE UNDERLINED.

Method	Publisher	Backbone	Shift Aware	SRCC	$R-\ell_2$
USDL [53]	CVPR'20	I3D	×	0.426	5.736
CoRe [8]	ICCV'21	I3D	×	0.471	5.402
TSA [48]	CVPR'22	I3D	X	0.452	5.533
CoRe-GOAT [3]	CVPR'23	I3D	×	0.494	5.072
HGCN [6]	TCSVT'23	I3D	X	0.471	4.954
CoFInAl [23]	IJCAI'24	I3D	1	0.552	4.586
PHI (Ours)	-	I3D	1	0.713	3.608
USDL [53]	CVPR'20	VST	×	0.473	5.076
CoRe [8]	ICCV'21	VST	Х	0.500	5.960
TSA [48]	CVPR'22	VST	X	0.475	4.778
CoRe-GOAT [3]	CVPR'23	VST	×	0.560	4.763
HGCN [6]	TCSVT'23	VST	X	0.671	6.564
CoFInAl [23]	IJCAI'24	VST	1	0.698	4.019
PHI (Ours)	-	VST	1	0.835	2.752

9

41

42

43

44

45

46

47

48

49

50

51

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

72

73

74

75

76

TABLE IV COMPUTATIONAL COMPARISON WITH EXISTING OPEN-SOURCE METHODS "-" DENOTES METHODS WITHOUT OFFLINE MODULES.

	Shift	FLOPs	Parame	eter (M)	Inference Time (ms)	
Method	Aware	(G)	Online	Offline		
ACTION-NET [15]	×	34.7500	28.08	-	305.2474	
GDLT [16]	×	0.1164	3.20	-	3.2249	
HGCN [6]	X	1.1201	0.50	-	6.7830	
CoFInAl [23]	✓	0.1178	3.70	-	3.8834	
PHI-half (Ours)	1	0.2637	2.80	4.60	5.6870	
PHI (Ours)	1	0.2637	3.00	4.60	5.6870	

of VST to better validate the generalizability of the proposed method on different backbones. This means that we did not fine-tune the parameters of the I3D backbone. As shown in Tables I and II, PHI is still superior to CoFInAI that has been well-tuned on the I3D backbone. As a result, it still has significant room for improvement through further fine-tuning, which could potentially lead to even better performance. This explains why the VST backbone shows greater improvement over I3D when compared to CoFInAI.

Although PHI is weaker than CoFInAl in some categories, PHI exhibits a 2.03% performance gain on the challenging Fis-V dataset, while achieving a modest 0.37% gain on the 52 simpler RG dataset. The smaller improvement observed on the RG dataset can be attributed to two key factors: its limited sample size and the relatively mild domain shift between the pre-training task and AQA. Specifically, RG comprises a smaller dataset compared to Fis-V (refer to the dataset details in Section V-A), which restricts the model's capacity to learn complex domain adaptation patterns. Furthermore, as illustrated in Figure 7(a), RG exhibits less pronounced feature misalignment, reducing the upper limit of performance improvement for extensive domain adaptation. In contrast, as can be seen in Figure 8(a), Fis-V demonstrates significant domain discrepancies, where PHI's domain shift mitigation mechanisms prove more impactful, resulting in a more substantial performance gain. Notably, PHI achieves a significant 42.9% reduction in $R-\ell_2$ error on RG, indicating PHI's robustness in effectively managing subtle feature discrepancies. Furthermore, the performance gap in CoFInAl can be attributed to the discretization of the continuous score space for classification. Our novel contribution focuses on addressing 71 feature-level discrepancies in AQA, which we have identified as crucial for achieving additional performance improvements, compared to the task-level one. By innovatively tackling these discrepancies, we have effectively enhanced the robustness and accuracy of AQA systems.

We further compare and analyze both solutions. While 77 CoFInAl and PHI addresses the domain shift in comple-78 mentary ways, they are fundamentally incompatible due to 79 their opposing alignment directions. As can be seen from 80 the different arrows in Figure 1(b), these two approaches 81 operate the alignment process in opposite directions, making 82 their direct combination infeasible. If pre-trained features are 83 already well-optimized and generalizable, task-level alignment 84 (CoFInAI) may be sufficient. However, our experimental re-85 sults in Tables I to III demonstrate that additional performance 86

and score types in all the three datasets, demonstrating its effectiveness in addressing the domain shift issue in long-term 2 AQA. Specifically, it excels in categories like Ball, Hoop, and 3 Ribbon on RG and across TES and PCS on Fis-V by a large 4 margin. For instance, PHI demonstrates notable performance 5 improvements, achieving a remarkable 9.65% and 6.14% cor-6 relation gain in the Ball and Hoop categories, respectively, on 7 the RG dataset compared to GDLT [16]. Overall, PHI delivers significant average correlation gains of 5.88%, 5.65%, and 24.44% on RG, Fis-V, and LOGO, respectively. These results 10 validate the benefit of employing the two-stage instruction to 11 address the domain shift issue. 12

The above statistics focus on unimodal comparisons and 13 exclude multi-modal methods. Below, we compare PHI with 14 recent multi-modal approaches. Notably, PAMFN [31] uses 15 only RGB data, while VATP-Net [51] leverages multi-modal 16 inputs. Despite being a unimodal method, PHI outperforms 17 PAMFN (0.711) by 13.9% and slightly surpasses the multi-18 modal method [47] on RG, as shown in Table I. On Fis-V, PHI 19 achieves an average SRCC 6.5% higher than PAMFN (0.755) 20 and 1.0% higher than VATP-Net (0.796), as shown in Table II. 21 These results highlight PHI's ability to handle domain shifts 22 without additional modalities, demonstrating that multi-modal 23 methods do not necessarily outperform unimodal methods 24 when domain shifts are effectively addressed. Future work 25 will explore extending PHI to multi-modal settings to further 26 enhance performance and generalization. The comparison with 27 various state-of-the-art methods positions PHI as a promising 28 solution for advancing AQA capabilities. 29

c) Comparison with Task-Level Discrepancies Solution: 30 This work categorizes domain shifts into two perspectives: 31 task-level discrepancies and feature-level discrepancies (see 32 Figure 1(b)). CoFInAl [23] primarily addressed task-level 33 discrepancies, whereas our work focuses on feature-level 34 discrepancies. Both approaches have demonstrated substan-35 tial performance improvements compared to shift-unaware 36 methods, as validated in Tables I to III, demonstrating the 37 effectivenes and necessity of explicit methods in mitigating 38 domain shifts in long-term AQA. In our experiments, we used 39 the optimal parameters for I3D from the well-tuned parameters 40

Sotting	Ball		Clubs		Ноор		Ribbon		Average	
Setting	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$	SRCC	$R-\ell_2$
PHI	0.818	2.187	0.803	2.149	0.812	2.199	0.805	2.744	0.810	2.300
PHI-half	$0.808 \downarrow 1.22\%$	2.027 ↓7.32%	0.790 ^{↓1.62%}	2.460 ^14.47%	0.788 ^{↓2.96%}	5.776 ^162.80%	$0.752 \downarrow 6.58\%$	2.900 ^{↑5.68%}	0.785 ^{↓3.09%}	3.291 ^{↑43.00%}
w/o GMF	0.802 ↓1.96%	1.942 ↓11.20%	0.784 ^{↓2.36%}	2.316 ^7.77%	0.789 ^{↓2.83%}	6.733 ^{206.27%}	0.756 \perp\$6.09\%	3.217 ^17.24%	0.783 43.33%	3.552 ^{↑54.43%}
w/o TESA	0.661 ↓19.20%	3.347 ^{↑53.09%}	0.637 \pressure 20.66%	3.904 ^{↑81.63%}	0.371 ↓54.32%	6.490 ^195.13%	0.550 \perp31.68%	5.016 ^{↑82.91%}	0.564 \perpart 30.37\%	4.689 ^103.87%
w/o LCR	$0.775 \downarrow 5.26\%$	3.105 ^42.00%	0.773 ^{↓3.74%}	$2.588^{\uparrow 20.40\%}$	0.789 ^{↓2.83%}	7.367 ^{↑235.09%}	$0.719 \downarrow 10.68\%$	3.671 ^33.81%	$0.765 \ ^{\downarrow 5.56\%}$	4.183 ^{↑81.87%}
w/o KL	0.806 ↓1.47%	1.991 ^{↓8.96%}	0.799 ↓0.50%	2.231 ^{↑3.81%}	$0.777 \downarrow 4.31\%$	7.005 ^218.56%	$0.778 \ ^{\downarrow 3.35\%}$	3.174 ^{↑15.68%}	0.790 ^{↓2.47%}	3.600 ^{↑56.52%}

TABLE V Ablation results on the RG dataset.

TABLE VI Ablation results on the Fis-V dataset.

Sotting	Sotting TES		Р	CS	Average		
Setting	SRCC	\mathbf{R} - ℓ_2	SRCC	\mathbf{R} - ℓ_2	SRCC	\mathbf{R} - ℓ_2	
PHI	0.726	2.543	0.867	1.656	0.804	2.178	
PHI-half	0.661 ^{↓9.07%}	3.050 19.88%	0.861 ↓1.15%	1.183 ^{↑28.58%}	0.780 ↓3.48%	2.117 ^{†2.79%}	
w/o GMF	0.668 ↓8.06%	3.413 ^34.33%	$0.857 \ ^{\downarrow 0.58\%}$	1.769 ^{↑6.72%}	0.780 ↓3.48%	2.591 ^{↑18.97%}	
w/o TESA	0.627 ↓13.87%	4.360 ^166.33%	$0.776 \downarrow 10.97\%$	3.714 ^123.67%	$0.709 \downarrow 11.84\%$	4.037 ^{↑84.81%}	
w/o LCR	$0.280 \downarrow 61.40\%$	5.576 ^{↑238.17%}	$0.294 \downarrow 66.16\%$	4.984 ^200.91%	$0.282 \downarrow 65.00\%$	5.280 ^142.05%	
w/o KL	$0.576 \downarrow 20.93\%$	2.437 ↓3.88%	$0.841 \downarrow 2.56\%$	2.564 ^{↑54.94%}	$0.780 \downarrow 3.48\%$	$2.501^{14.74\%}$	

gains are achievable through PHI, indicating that pre-trained features still require adaptation to better suit AQA. In the 2 future, we aim to explore the potential of combining both з approaches to fully leverage their respective strengths and 4 achieve even higher performance. Moreover, PHI is designed 5 as a modular enhancement and can be integrated into other 6 baseline methods that do not explicitly consider domain shift. By refining feature alignment, PHI enhances model robustness 8 and generalization, making it a versatile tool for improving 9 AQA performance across various architectures. 10

d) Computational Efficiency and Model Complexity: To 11 provide a more comprehensive comparison, we evaluated PHI 12 and existing long-term AQA methods under identical condi-13 tions, with a focus on computational efficiency. The results are 14 summarized in Table IV. PHI consists of online and offline 15 components, strategically designed to balance computational 16 efficiency with assessment accuracy. The online module (flow 17 path) handles real-time inference with minimal latency, while 18 the offline module (TETE) refines feature representations, thus 19 reducing the computational burden during online inference. 20 Combining the results from Tables I to III, the distillation 21 design enables it to achieve a competitive balance between 22 performance and efficiency compared to state-of-the-art mod-23 els. 24

Notably, CoFInAl achieves the lowest FLOPs (0.1178G) 25 and inference time (3.8834ms) among shift-aware methods, 26 underscoring its computational efficiency. However, PHI com-27 pensates for a slightly higher computational cost by leveraging 28 a progressive instruction strategy. This approach improves 29 assessment accuracy with only a 0.1470G increase in FLOPs 30 and a 2.4621ms additional inference delay compared to GDLT. 31 Despite this, PHI outperforms in terms of online parame-32 ters and quality assessment performance, while maintaining 33 a reasonable computational footprint. PHI introduces only 34 a slight increase in FLOPs and parameters, while signifi-35 cantly improving the assessment accuracy. PHI reduces FLOPs 36 (by 99.2%), parameters (by 89.3%), and inference time (by 37

98.1%) compared to ACTION-NET, and also achieves lower 38 FLOPs and competitive inference time relative to HGCN. 39 Additionally, PHI benefits from an offline distillation mech-40 anism, allowing the TETE module to transfer knowledge to 41 a lightweight flow model, further reducing online parameters 42 and computational cost. In summary, while CoFInAl excels 43 in computational efficiency, PHI strikes an optimal balance 44 between computational cost and assessment accuracy, making 45 it a practical solution for practical AQA applications. 46

2) Ablation Studies: This study aims to investigate the individual contributions of the core components of PHI, particularly focusing on GMF (see Hypothesis 1), TESA (see Section IV-A2), LCR (see Hypothesis 2), and the use of KL divergence loss (see Equation (15)). Tables V and VI and Figure 5 present the results on RG and Fis-V.

a) Validation of Hypothesis 1: To validate the effective-53 ness of Hypothesis 1, we compare removing GMF (*w/o GMF*) 54 and TETE (w/o TETE), respectively. On the one hand, when 55 removing GMF while retaining only TETE (w/o GMF), we 56 observe a decrease in performance across all categories, with 57 an average SRCC decrease of 3.33% and an average $R-\ell_2$ 58 increase of 54.43% on the RG dataset (see Table V) and 59 an average SRCC decrease of 3.48% and an average $R-\ell_2$ 60 increase of 18.97% on the Fis-V dataset (see Table VI). These 61 results indicate that the desired features estimated by TETE 62 alone contain inaccuracies due to the inherent domain gap. The 63 GMF module plays a crucial role in refining these estimations 64 by progressively reducing the domain shift, thereby enhancing 65 the reliability of feature representations and improving action 66 assessment performance. On the one hand, by replacing TESA 67 with vanilla attention (w/o TESA), we observe a more substan-68 tial decrease in performance, with an average SRCC decrease 69 of 30.37% on the RG dataset (see Table V) and an average 70 SRCC decrease of 11.84% on the Fis-V dataset (see Table VI). 71 This emphasizes the critical role of TESA in capturing long-72 range dependencies efficiently, which is essential for accurate 73 AQA. These results collectively validate the effectiveness of 74

47

48

49

50

51



Fig. 5: Results of (a) SRCC and (b) $R-\ell_2$ on the impact of different steps. The symbol " \uparrow " indicates higher is better, while the symbol " \downarrow " indicates lower is better.

Hypothesis 1

b) Validation of Hypothesis 2: To validate the effective-2 ness of Hypothesis 2, we compare removing LCR (w/o LCR) 3 and KL (w/o KL), respectively. On the one hand, removing 4 LCR (w/o LCR) results in a notable decrease in performance, 5 with an average SRCC decrease of 5.56% and an average 6 $R-\ell_2$ increase of 81.87% on the RG dataset (see Table V) and an average SRCC decrease of 65.00% and an average R-8 ℓ_2 increase of 142.15% on the Fis-V dataset (see Table VI). g This demonstrates that LCR plays a crucial role in learning 10 representations focusing on fine-grained cues, which are vital 11 for mitigating domain shift and improving AQA performance. 12 On the one hand, replacing the KL divergence loss with MSE 13 (w/o KL) leads to a significant decrease in performance, with 14 an average R- ℓ_2 increase of 56.52% on the RG dataset (see 15 Table V) and an average $R-\ell_2$ increase of 14.74% on the Fis-V 16 dataset (see Table VI). This highlights the effectiveness of the 17 KL divergence loss in guiding the model to learn more robust 18 representations aligned with AQA. These results collectively 19 validate the effectiveness of Hypothesis 2. 20

c) Impact of the Model Size: As shown in Tables V 21 and VI, reducing the parameter size of the flow network ϕ by 22 half (PHI-half) results in a noticeable decrease in performance 23 across all categories. Specifically, we observe an average 24 decrease of 3.09% in SRCC and an average increase of 43.00% 25 in R- ℓ_2 on the RG dataset and an average decrease of 3.48% 26 in SRCC and an average increase of 2.79% in R- ℓ_2 on the 27 RG dataset. Combined with the reported results in Tables I 28 and II, we observe that PHI-half still outperforms some strong 29 baselines [16], [6]. This finding shows the importance of 30 model size in maintaining performance levels, suggesting that 31 a larger parameter space with simple MLPs contributes to 32 better overall performance, indicating the effectiveness of 33 PHI's network design. 34

d) Impact of Different Steps: The number of flow steps 35 plays a crucial role in refining pre-trained features for AQA. 36 Adjusting this parameter influences the model's ability to 37 reduce the domain gap between pre-trained features and AQA 38 tasks (see Figure 5). Firstly, while increasing the number of 39 steps yields slight improvements, the gains are not always sub-40 stantial. To provide an intuitive understanding, Figure 6 illus-41 trates the conceptual differences between one-step and multi-42 step flows. The one-step approach (see Figure 6(a)) achieves 43 competitive performance due to its direct alignment but is 44 more susceptible to deviations. In contrast, the multi-step flow 45 (see Figure 6(b)) provides more controlled and gradual refine-46



Fig. 6: Conceptual illustrations of (a) one-step and (b,c) multi-step flows. In (a), the one-step flow allows for a fast and direct reduction of domain gaps, making it computationally efficient. However, it may lead to larger errors. Suppose the second state deviates from the desired direction with the same degree θ . In contrast, the multistep flow (b) provides a more robust alignment through gradual refinements, reducing the risk of large errors. However, in certain cases (c), multi-step alignment can accumulate errors, potentially leading to worse performance than the one-step approach.

 TABLE VII

 Results of different training strategies.

Strategy	RG	Fis-V
Two-stage	0.810	0.804
One-stage	$0.807 \ ^{\downarrow 0.37\%}$	$0.800 \ ^{\downarrow 0.50\%}$

ment, potentially reducing large alignment errors. However, 47 multi-step alignment can also introduce accumulative errors 48 (see Figure 6(c)), especially when the initial step already aligns 49 features effectively. This explains why, in some cases, multi-50 step flows offer only marginal improvements or even perform 51 worse than the one-step approach (in Figure 5). Nonethe-52 less, multi-step alignment enhances stability and robustness, 53 particularly in scenarios where direct alignment might lead 54 to suboptimal feature adaptation. Additionally, increasing the 55 number of steps allows for a more gradual transformation of 56 initial features into task-specific representations, potentially 57 improving accuracy and reliability. However, this comes at the 58 cost of higher computational complexity and longer training 59 time. Conversely, reducing the number of steps accelerates 60 training but may limit the model's ability to effectively align 61 with AQA tasks. Importantly, our method is designed to 62 support both one-step and multi-step alignment, offering flexi-63 bility depending on the computational constraints and accuracy 64 requirements. The absence of a clear performance trend across 65 different step settings further highlights the robustness and 66 adaptability of our approach. 67

e) Impact of Different Training Strategies: We adopt a two-stage training process in our approach. Initially, we train the TETE to refine the estimation of desired features, focusing on obtaining accurate representations essential for the flow model. Subsequently, we integrate these refined features into our overall framework, enabling joint training with components like GMF. To demonstrate the efficacy of the two-stage training approach, we conducted experiments comparing it to a one-stage joint training process. The results, shown in Table VII, reveal that the two-stage process outperforms the one-stage process, with improvements of 0.37% on the RG dataset and 0.50% on the Fis-V dataset. This highlights the clear advantages of the two-stage training strategy.

3) Qualitative and Quantitative Results: We present a diverse range of visualizations to demonstrate both the qualitative and quantitative performance of our method.

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

a) Visualization of the Domain Shift: We visually demonstrate the effectiveness of our PHI method in mitigating 2 domain shifts through feature distribution visualizations in 3 the latent space and correlation analyses between predicted 4 and ground truth scores. Specifically, we employ the t-SNE toolbox [54] to generate feature distribution plots, which are 6 presented in the first rows of Figures 7 and 8. Additionally, we illustrate the correlation between predicted and actual 8 scores through scatter plots, as shown in the second rows of Figures 7 and 8. To assess the model's generalization 10 capability, we visualize two action categories, Hoop and PCS, 11 selected from the RG and PCS datasets, respectively. In the 12 feature distribution plots, different score ranges (or grades) 13 are delineated using an SVC classifier. Due to variations in 14 dataset scales, we categorize the test samples into four score 15 ranges for Hoop (RG), assigning labels from 0 to 3, and 16 six score ranges for PCS (Fis-V), assigning labels from 0 to 17 5. Improved feature clustering, where samples of the same 18 grade (represented by similar color shading) occupy the same 19 region, indicates a more structured feature space for action 20 assessment. Furthermore, we provide comparisons with GDLT 21 [16] and CoFInAl [23] to highlight the advantages of our 22 approach. Since the Fis-V dataset contains a larger number of 23 test samples compared to RG, it provides a more robust and 24 reliable evaluation of model performance and generalizability. 25 Therefore, we primarily focus on comparisons conducted on 26 the PCS (Fis-V) dataset for a more comprehensive assessment 27 of our method in the following. 28

Specifically, the features extracted by the VST backbone, as 29 shown in Figure 8(a), exhibit a mixed distribution, indicating 30 difficulties in distinguishing between different score ranges. 31 This suggests that the pre-trained backbone may not be well-32 suited for the AQA task, leading to significant domain shift 33 issues. In Figure 8(b), the feature distribution of GDLT appears 34 confused, reflecting ineffective feature learning. In contrast, 35 Figures 8(c) and 8(d) illustrate the feature distributions of 36 CoFInAl and PHI, respectively, which display more distinct 37 clustering. This clearer separation facilitates the identification 38 of samples within each score range, indicating that both 39 CoFInAl and PHI effectively mitigate domain shift. While the 40 advantage of PHI can be observed by comparing the feature 41 plots in Figures 7(c) and 7(d), the t-SNE visualizations in 42 Figures 8(c) and 8(d) do not provide a definitive comparison 43 between the domain shift mitigation capabilities of PHI and 44 CoFInAl. Notably, CoFInAl suffers from a loss of precision 45 in score prediction, which affects its fine-grained assessment 46 capability. In contrast, PHI maintains higher precision, leading 47 to improved reliability in action assessment. To further clarify 48 this, we demonstrate that PHI outperforms CoFInAl in the 49 following analysis. 50

Figures 8(e), 8(f) and 8(h) compare correlation plots be-51 tween GDLT, CoFInAl, and PHI. In Figure 8(e), the correlation 52 line $(\hat{s}_i = 0.28s_i + 18.72)$ of GDLT shows a deviation from the 53 ideal correlation line $(\hat{s}_i = s_i)$, indicating a weak correlation 54 with ground truth scores. In Figure 8(f), the correlation line 55 $(\hat{s}_i = 0.34s_i + 16.95)$ of CoFInAl shows a smaller deviation 56 from the ideal line, suggesting a stronger correlation with 57 ground truth scores compared to GDLT. Notably, Figure 8(h) 58



Fig. 7: Visualization depicting the mitigation of domain shift on the Hoop (RG) dataset: t-SNE feature distribution plots (a, b, c, d) and correlation comparison plots (e, f, g) of GDLT [16], CoFInAl [23], and our PHI method. The dataset is split into four grades.



Fig. 8: Visualization depicting the mitigation of domain shift on the PCS (Fis-V) dataset: t-SNE feature distribution plots (a, b, c, d) and correlation comparison plots (e, f, g) of GDLT [16], CoFInAl [23], and our PHI method. The dataset is split into six grades.

shows that the correlation line of PHI ($\hat{s}_i = 0.46s_i + 12.80$) is the closest to the ideal line, demonstrating the highest correlation with ground truth scores among the three methods. This suggests that PHI achieves a higher correlation with ground truth scores compared to GDLT and CoFInAl, further highlighting its superiority in mitigating domain shift and improving AQA performance. Overall, the visualizations in Figure 8 provide compelling evidence of PHI's effectiveness in mitigating domain shift and improving long-range AQA performance compared to existing methods.

b) Visualization of Attention Weights: To gain deeper insights into the attention mechanism within the TETE module, we visualize the attention weights and the highlighted clips for two representative samples from the RG and Fis-V datasets. Figures 9(a) and 9(b) illustrate the attention weight distributions for each clip in the Club (RG) and PCS (Fis-V) models, respectively. The first row in each subfigure presents the normalized attention weights assigned to all action clips, revealing the varying levels of importance attributed to different segments of the action. The following rows highlight the top three clips that received the highest attention scores, which are crucial for the model's decision-making. In Figure 9(a), the model predicts a score of 14.30, while the ground-truth score is 14.52, resulting in a minimal error of 0.22. Similarly, in Figure 9(b), the model predicts a score of 27.78, compared to the ground-truth score of 26.68, yielding an error of 1.10. For example, in Figure 9(b), the three most attended clips are 39,

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84





(b) Sample #056 from the Fis-V (PCS) dataset

Fig. 9: Visualization of attention weights and highlighted clips for samples on (a) the RG (Club) dataset and (b) the Fis-V (PCS) dataset.



Fig. 10: Visualization of the distance matrix on the Ball (RG) dataset.

52, and 53, corresponding to key subactions where the player
executes precise hand-leg coordination. These visualizations
provide valuable insights into how the AQA model prioritizes
different temporal segments within an action sequence. By
identifying the most critical moments, we enhance our understanding of the model's interpretability and its evaluation
process for action quality assessment.

c) Visualization of the Distance Matrix: In Figure 10, 8 the heatmap of the distance matrix is depicted. The highest 9 value in the matrix, such as (12, 27), signifies the greatest 10 distance between the 12-th and 27-th actions. To further assess 11 the efficacy, the two actions and their respective heatmaps are 12 visualized. These heatmaps in Figures 10(c) and 10(d) provide 13 insights into the domain shift issue. The minimal difference 14 between the predicted and ground-truth scores indicates the 15 effectiveness of our method. For instance, with a ground-truth 16 score of $s_{12} = 8.40$ and the predicted score is $\hat{s}_{12} = 8.64$. 17 Additionally, observing the score distance between the two 18 19 videos reveals a wide range, consistent with their feature distance, demonstrating the effectiveness of our PHI method. 20

21 VI. CONCLUSIONS AND DISCUSSIONS

In this work, we identify and analyze task-level and featurelevel domain shifts in long-term AQA and propose PHI as a hierarchical adaptation framework to address feature-level dis-24 crepancies. Rather than a mere extension of existing methods, 25 PHI introduces a novel integration of shallow-to-deep adapta-26 tion and coarse-to-fine alignment strategies to enhance AQA 27 performance. The shallow-to-deep adaptation strategy, enabled 28 by GMF, effectively reduces domain gaps while maintain-29 ing computational efficiency. Simultaneously, the coarse-to-30 fine alignment mechanism, facilitated by LCR, refines coarse 31 features extracted from pre-trained models, aligning them with 32 fine-grained representations crucial for AQA. Experimental 33 results on three representative long-term AQA datasets demon-34 strate the significant effectiveness of PHI, underscoring the im-35 portance of mitigating feature-level discrepancies in improving 36 AQA performance. Notably, compared to the task-level adap-37 tation method CoFInAI, PHI exhibits superior performance 38 in mitigating domain shifts and enhancing AQA accuracy, 39 emphasizing the critical role of feature-level alignment in 40 long-term AQA tasks. Furthermore, the hierarchical adaptation 41 framework of PHI is highly generalizable beyond AQA, with 42 potential applications in rehabilitation analysis, sports motion 43 scoring, and movement disorder diagnosis. The principles of 44 shallow-to-deep adaptation and coarse-to-fine alignment also 45 extend to domains such as multi-modal alignment. By address-46 ing domain shifts through hierarchical feature refinement, PHI 47 provides a theoretically grounded and versatile framework, 48 paving the way for future research across multiple fields.

Despite its strong performance, PHI has several limitations, 2 each suggesting directions for future work. First, the auto-3 regressive nature of GMF introduces cumulative errors, which may progressively degrade prediction accuracy over multiple steps. Future research will focus on advanced optimization 6 strategies and novel regularization techniques to mitigate these challenges. Second, while our clip alignment method effectively minimizes discrepancies between actions, it may lead to information loss in cases of significant temporal variation. 10 Future work will explore adaptive alignment strategies to 11 enhance temporal correlation capture while maintaining com-12 putational efficiency. Third, although PHI and CoFInAI ad-13 dress domain shifts from complementary perspectives, feature-14 level and task-level alignment, respectively, their integration 15 into a unified framework remains an open challenge. Future 16 research will investigate synergistic strategies to combine 17 these approaches for a more comprehensive domain adaptation 18 solution. Additionally, PHI is currently designed for unimodal 19 inputs, limiting its applicability in multi-modal scenarios. Ex-20 tending the framework to support multi-modal integration will 21 be a key avenue for future work, broadening its utility across 22 diverse applications. Finally, PHI is tailored for mitigating the 23 domain shift issue in long-term AQA tasks, which may not 24 be directly applicable to short-term AOA scenarios. Future 25 research will explore adaptations to extend its applicability to 26 short-term AOA, further increasing its versatility. Addressing 27 these limitations will enhance the robustness, generalizability, 28 and applicability of PHI, advancing the field of action quality 29

30 assessment and its related domains.

31

36

37

38

39

REFERENCES

- [1] K. Zhou, L. Wang, X. Zhang, H. P. Shum, F. W. Li, J. Li, and X. Liang,
 "Magr: Manifold-aligned graph regularization for continual action quality assessment," in *Proceedings of the 18th European Conference on Computer Vision*, 2024.
 - [2] Y. Ji, L. Ye, H. Huang, L. Mao, Y. Zhou, and L. Gao, "Localizationassisted uncertainty score disentanglement network for action quality assessment," in *Proceedings of the 31st ACM International Conference* on Multimedia, pp. 8590–8597, 2023.
- [3] S. Zhang, W. Dai, S. Wang, X. Shen, J. Lu, J. Zhou, and Y. Tang,
 "Logo: A long-form video dataset for group action quality assessment,"
 in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2405–2414, 2023.
- [4] J. Xu, S. Yin, G. Zhao, Z. Wang, and Y. Peng, "Fineparser: A finegrained spatio-temporal action parser for human-centric action quality assessment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14628–14637, 2024.
- [5] Y.-M. Li, L.-A. Zeng, J.-K. Meng, and W.-S. Zheng, "Continual action assessment via task-consistent score-discriminative feature distribution modeling," *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.
- [6] K. Zhou, Y. Ma, H. P. Shum, and X. Liang, "Hierarchical graph convolutional networks for action quality assessment," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 12, pp. 7749–7763, 2023.
- [7] M. Li, H.-B. Zhang, Q. Lei, Z. Fan, J. Liu, and J.-X. Du, "Pairwise contrastive learning network for action quality assessment," in *European Conference on Computer Vision*, pp. 457–473, Springer, 2022.
- [8] X. Yu, Y. Rao, W. Zhao, J. Lu, and J. Zhou, "Group-aware contrastive re gression for action quality assessment," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 7919–7928, 2021.
- [9] L. Yao, Q. Lei, H. Zhang, J. Du, and S. Gao, "A contrastive learning network for performance metric and assessment of physical rehabilitation exercises," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 2023.

- [10] K. Zhou, R. Cai, Y. Ma, Q. Tan, X. Wang, J. Li, H. P. Shum, F. W. Li, S. Jin, and X. Liang, "A video-based augmented reality system for human-in-the-loop muscle strength assessment of juvenile dermatomyositis," *IEEE Transactions on Visualization and Computer Graphics*, vol. 29, no. 5, pp. 2456–2466, 2023.
- [11] S. Deb, M. F. Islam, S. Rahman, and S. Rahman, "Graph convolutional networks for assessment of physical rehabilitation exercises," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 30, pp. 410–419, 2022.
- [12] C. K. Ingwersen, A. Xarles, A. Clapés, M. Madadi, J. N. Jensen, M. R. Hannemose, A. B. Dahl, and S. Escalera, "Video-based skill assessment for golf: Estimating golf handicap," in *Proceedings of the* 6th International Workshop on Multimedia Content Analysis in Sports, pp. 31–39, 2023.
- [13] L. Trinh, T. Chu, Z. Cui, A. Malpani, C. Yang, I. Dalieh, A. Hui, O. Gomez, Y. Liu, and A. Hung, "Self-supervised sim-to-real kinematics reconstruction for video-based assessment of intraoperative suturing skills," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 708–717, Springer, 2023.
- [14] X. Ding, X. Xu, and X. Li, "Sedskill: Surgical events driven method for skill assessment from thoracoscopic surgical videos," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 35–45, Springer, 2023.
- [15] L.-A. Zeng, F.-T. Hong, W.-S. Zheng, Q.-Z. Yu, W. Zeng, Y.-W. Wang, and J.-H. Lai, "Hybrid dynamic-static context-aware attention network for action assessment in long videos," in *Proceedings of the 28th ACM international conference on multimedia*, pp. 2526–2534, 2020.
- [16] A. Xu, L.-A. Zeng, and W.-S. Zheng, "Likert scoring with grade decoupling for long-term action assessment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3232–3241, 2022.
- [17] A. Dadashzadeh, S. Duan, A. Whone, and M. Mirmehdi, "Pecop: Parameter efficient continual pretraining for action quality assessment," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 42–52, 2024.
- [18] P. Parmar and B. Morris, "Action quality assessment across multiple actions," in 2019 IEEE winter conference on applications of computer vision (WACV), pp. 1468–1476, IEEE, 2019.
- [19] Y. Bai, D. Zhou, S. Zhang, J. Wang, E. Ding, Y. Guan, Y. Long, and J. Wang, "Action quality assessment with temporal parsing transformer," in *European conference on computer vision*, pp. 422–438, Springer, 2022.
- [20] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, pp. 6299–6308, 2017.
- [21] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," in *Proceedings of the IEEE/CVF conference on computer* vision and pattern recognition, pp. 3202–3211, 2022.
- [22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [23] K. Zhou, J. Li, R. Cai, L. Wang, X. Zhang, and X. Liang, "Cofinal: Enhancing action quality assessment with coarse-to-fine instruction alignment," in *Proceedings of the 33rd International Joint Conference* on Artificial Intelligence, 2024.
- [24] X. Liu, C. Gong, and Q. Liu, "Flow straight and fast: Learning to generate and transfer data with rectified flow," *arXiv preprint arXiv:2209.03003*, 2022.
- [25] Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, "Flow matching for generative modeling," *arXiv preprint arXiv:2210.02747*, 2022.
- [26] P. Parmar and B. Tran Morris, "Learning to score olympic events," in Proceedings of the IEEE conference on computer vision and pattern recognition workshops, pp. 20–28, 2017.
- [27] H. Pirsiavash, C. Vondrick, and A. Torralba, "Assessing the quality of actions," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI 13*, pp. 556–571, Springer, 2014.
- [28] S. Wang, D. Yang, P. Zhai, C. Chen, and L. Zhang, "Tsa-net: Tube self-attention network for action quality assessment," in *Proceedings of* the 29th ACM international conference on multimedia, pp. 4902–4910, 2021.
- [29] C. Xu, Y. Fu, B. Zhang, Z. Chen, Y.-G. Jiang, and X. Xue, "Learning to score figure skating sport videos," *IEEE transactions on circuits and systems for video technology*, vol. 30, no. 12, pp. 4578–4590, 2019.
- [30] X. Wang, J. Li, and H. Hu, "Skeleton-based action quality assessment via partially connected lstm with triplet losses," in *Chinese Conference*

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

2

3

4

7

22

on Pattern Recognition and Computer Vision (PRCV), pp. 220-232, Springer, 2022.

- [31] L.-A. Zeng and W.-S. Zheng, "Multimodal action quality assessment," IEEE Transactions on Image Processing, vol. 33, pp. 1600–1613, 2024.
- [32] S. Wang, D. Yang, P. Zhai, Q. Yu, T. Suo, Z. Sun, K. Li, and 5 6 L. Zhang, "A survey of video-based action quality assessment," in 2021 International Conference on Networking Systems of AI (INSAI), pp. 1-9, IEEE, 2021. 8
- [33] P. Parmar and B. T. Morris, "What and how well you performed? a mul-9 titask learning approach to action quality assessment," in Proceedings of 10 11 the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 304-313, 2019. 12
- [34] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning 13 14 spatiotemporal features with 3d convolutional networks," in Proceedings of the IEEE international conference on computer vision, pp. 4489-15 4497 2015 16
- [35] J.-H. Pan, J. Gao, and W.-S. Zheng, "Action assessment by joint relation 17 graphs," in Proceedings of the IEEE/CVF international conference on 18 19 computer vision, pp. 6331-6340, 2019.
- Z. Du, D. He, X. Wang, and Q. Wang, "Learning semantics-guided repre-20 [36] sentations for scoring figure skating," IEEE Transactions on Multimedia, 21 vol. 26, pp. 4987-4997, 2024.
- [37] J. Xia, M. Zhuge, T. Geng, S. Fan, Y. Wei, Z. He, and F. Zheng, 23 24 "Skating-mixer: Long-term sport audio-visual modeling with mlps," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 37, 25 pp. 2901–2909, 2023. 26
- 27 [38] A. Tong, N. Malkin, G. Huguet, Y. Zhang, J. Rector-Brooks, K. Fatras, G. Wolf, and Y. Bengio, "Improving and generalizing flow-based 28 generative models with minibatch optimal transport," arXiv preprint 29 arXiv:2302.00482, 2023. 30
- [39] M. S. Albergo and E. Vanden-Eijnden, "Building normalizing flows with 31 32 stochastic interpolants," arXiv preprint arXiv:2209.15571, 2022.
- [40] D. Onken, S. W. Fung, X. Li, and L. Ruthotto, "Ot-flow: Fast and accu-33 34 rate continuous normalizing flows via optimal transport," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 35, pp. 9223-35 9232, 2021. 36
- 37 [41] I. Kobyzev, S. J. Prince, and M. A. Brubaker, "Normalizing flows: An introduction and review of current methods," IEEE transactions on 38 pattern analysis and machine intelligence, vol. 43, no. 11, pp. 3964-39 40 3979, 2020.
- [42] K. Neklyudov, D. Severo, and A. Makhzani, "Action matching: A 41 42 variational method for learning stochastic dynamics from samples," arXiv preprint arXiv:2210.06662, 2022. 43
- [43] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," 44 45 Advances in neural information processing systems, vol. 33, pp. 6840-6851 2020 46
- [44] Z. Chang, G. A. Koulieris, and H. P. Shum, "On the design fundamentals 47 of diffusion models: A survey," arXiv preprint arXiv:2306.04542, 2023. 48
- 49 [45] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and 50 B. Poole, "Score-based generative modeling through stochastic differential equations," arXiv preprint arXiv:2011.13456, 2020. 51
- 52 [46] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," arXiv preprint arXiv:2010.02502, 2020. 53
- K. Gedamu, Y. Ji, Y. Yang, J. Shao, and H. T. Shen, "Fine-grained [47] 54 spatio-temporal parsing network for action quality assessment," IEEE 55 Transactions on Image Processing, vol. 32, pp. 6386-6400, 2023. 56
- 57 [48] J. Xu, Y. Rao, X. Yu, G. Chen, J. Zhou, and J. Lu, "Finediving: A fine-grained dataset for procedure-aware action quality assessment," 58 in Proceedings of the IEEE/CVF conference on computer vision and 59 60 pattern recognition, pp. 2949-2958, 2022.
- [49] F. Raiber and O. Kurland, "Kullback-leibler divergence revisited," in 61 Proceedings of the ACM SIGIR international conference on theory of 62 information retrieval, pp. 117-124, 2017. 63
- [50] X. Ke, H. Xu, X. Lin, and W. Guo, "Two-path target-aware contrastive 64 65 regression for action quality assessment," Information Sciences, vol. 664, p. 120347, 2024. 66
- [51] K. Gedamu, Y. Ji, Y. Yang, J. Shao, and H. T. Shen, "Visual-semantic 67 alignment temporal parsing for action quality assessment," IEEE Trans-68 actions on Circuits and Systems for Video Technology, 2024. 69
- K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image 70 [52] recognition," in CVPR, pp. 770-778, 2016. 71
- Y. Tang, Z. Ni, J. Zhou, D. Zhang, J. Lu, Y. Wu, and J. Zhou, 72 [53] 73 "Uncertainty-aware score distribution learning for action quality assessment," in Proceedings of the IEEE/CVF conference on computer vision 74 and pattern recognition, pp. 9839-9848, 2020. 75
- [54] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne.," 76 77 Journal of machine learning research, vol. 9, no. 11, 2008.





Kanglei Zhou is a Ph.D. candidate in the School of Computer Science and Engineering at Beihang University, specializing in action quality assessment and augmented reality. From February to August 2024, he was a visiting student in the Department of Computer Science at Durham University. He received his Bachelor's degree in the College of Computer and Information Engineering from Henan Normal University in 2020.

Hubert P. H. Shum (Senior Member, IEEE) is a Professor of Visual Computing and the Director of Research of the Department of Computer Science at Durham University, specialising in modelling spatiotemporal information with responsible AI. He is also a Co-Founder and the Co-Director of Durham University Space Research Centre. Before this, he was an Associate Professor/Senior Lecturer at Northumbria University and a Postdoctoral Researcher at RIKEN Japan. He received his PhD degree from the University of Edinburgh. He chaired conferences

such as Pacific Graphics, BMVC and SCA, and has authored over 180 research publications.





ics, in 2020.



Frederick W. B. Li received a B.A. and an M.Phil. degree from Hong Kong Polytechnic University, and a Ph.D. degree from the City University of Hong Kong. He is currently an Associate Professor at Durham University, researching computer graphics, deep learning, collaborative virtual environments, and educational technologies. He is also an Associate Editor of Frontiers in Education and an Editorial Board Member of Virtual Reality & Intelligent Hardware. He chaired conferences such as ISVC and ICWL.

Xingxing Zhang received the BE and PhD degrees from the Institute of Information Science, Beijing Jiaotong University, in 2015 and 2020, respectively. She was also a visiting student with the Department of Computer Science, University of Rochester, from 2018 to 2019. She was a postdoc with the Department of Computer Science and Technology, Tsinghua University, from 2020 to 2022. Her research interests include continual learning and zero/fewshot learning. She has received the excellent PhD thesis award from the Chinese Institute of Electron-

Xiaohui Liang received his Ph.D. degree in computer science and engineering from Beihang University, China. He is currently a Professor, working in the School of Computer Science and Engineering at Beihang University. His main research interests include computer graphics and animation, visualization, and virtual reality.

131

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129



Citation on deposit: Zhou, K., Shum, H. P. H., Li, F. W. B., Zhang, X., & Liang, X. (in press). PHI: Bridging Domain Shift in Long-Term Action Quality Assessment via Progressive Hierarchical Instruction. IEEE Transactions on Image Processing

For final citation and metadata, visit Durham Research Online URL:

https://durham-repository.worktribe.com/output/3964025

Copyright statement: This accepted manuscript is licensed under the Creative Commons Attribution 4.0 licence. https://creativecommons.org/licenses/by/4.0/