

ContentCTR: Frame-level Live Streaming Click-Through Rate Prediction with Multimodal Transformer

Jiaxin Deng*
National Laboratory of Pattern
Recognition, Institute of Automation,
Chinese Academy of Sciences
Haidian Qu, Beijing Shi, China
dengjiaxin2022@ia.ac.cn

Dong Shen
KuaiShou Inc.
Haidian Qu, Beijing Shi, China
shendong@kuaishou.com

Shiyao Wang
KuaiShou Inc.
Haidian Qu, Beijing Shi, China
wangshiyao08@kuaishou.com

Xiangyu Wu
KuaiShou Inc.
Haidian Qu, Beijing Shi, China
wuxiangyu@kuaishou.com

Fan Yang
KuaiShou Inc.
Haidian Qu, Beijing Shi, China
yangfan@kuaishou.com

Guorui Zhou
KuaiShou Inc.
Haidian Qu, Beijing Shi, China
zhouguorui@kuaishou.com

Gaofeng Meng[†]
National Laboratory of Pattern
Recognition, Institute of Automation,
Chinese Academy of Sciences
Haidian Qu, Beijing Shi, China
gfmeng@nlpr.ia.ac.cn

ABSTRACT

In recent years, live streaming platforms have gained immense popularity as they allow users to broadcast their videos and interact in real-time with hosts and peers. Due to the dynamic changes of live content, accurate recommendation models are crucial for enhancing user experience. However, most previous works treat the live as a whole item and explore the Click-through-Rate (CTR) prediction framework on *item-level*, neglecting that the dynamic changes that occur even within the same live room. In this paper, we proposed a ContentCTR model that leverages multimodal transformer for *frame-level* CTR prediction. First, we present an end-to-end framework that can make full use of multimodal information, including visual frames, audio, and comments, to identify the most attractive live frames. Second, to prevent the model from collapsing into a mediocre solution, a novel pairwise loss function with first-order difference constraints is proposed to utilize the contrastive information existing in the highlight and non-highlight frames. Additionally, we design a temporal text-video alignment module based on Dynamic Time Warping to eliminate noise caused by the ambiguity and non-sequential alignment of visual and textual information. We conduct extensive experiments on both real-world

scenarios and public datasets, and our ContentCTR model outperforms traditional recommendation models in capturing real-time content changes. Moreover, we deploy the proposed method on our company platform, and the results of online A/B testing further validate its practical significance.

CCS CONCEPTS

• Information systems → Data mining.

KEYWORDS

stream highlight detection, multi-modal learning, click-through rate prediction

ACM Reference Format:

Jiaxin Deng, Dong Shen, Shiyao Wang, Xiangyu Wu, Fan Yang, Guorui Zhou, and Gaofeng Meng. 2023. ContentCTR: Frame-level Live Streaming Click-Through Rate Prediction with Multimodal Transformer. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Live streaming platform represent a new type of online interaction and experienced rapid growth in recent years. On such platform, more and more users are attracted to live streaming for entertainment, while streamers receive rewards from both the audiences and the platform. This new form of interaction and entertainment has motivated researchers to study emerging issues such as gift-sending mechanisms [38], E-Commerce events [34], and other live streaming practices. Due to the large number of streamers and the constantly changing content, an accurate recommendation algorithm is crucial for enhancing user experience. [35, 36] focus on capturing mutual information between streamers and viewers to obtain better representations. [25] propose a self-attentive model

*Interns at MMU, KuaiShou Inc.

[†]Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2023 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

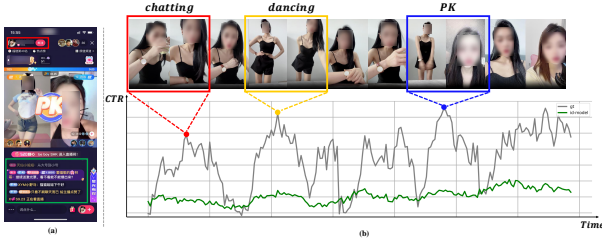


Figure 1: (a) shows the screenshots of the platform mobile app. There exist four kinds of modality: the visual feature of streaming frames, the speech from streamer, the ID embedding of streamer (red box) and the bullet comments from the user (green box). (b) The grey line shows the ground truth CTR changes of a dancing streamer while the green line represents the predicted CTR of ID-based recommendation model. And three distinct highlight moments are marked by red, yellow, and blue points.

to address the dynamic availability and repeat consumption of live streaming. [34] present a novel Live Stream E-Commerce Graph Neural Network framework to learn the tripartite interaction among streamers, users and products.

The broadcast of real-time content implies the live is constantly changing, and the performance of various segments may differ significantly. Although most previous work [25, 34–36] treat live broadcasts as items and model them at the item-level, we believe that understanding the content of individual frames and making more refined, frame-level estimates is crucial. Additionally, identifying the most captivating frames from the stream and presenting them to potentially interested users can considerably enhance user engagement, attract more users, and increase revenue for live-streaming platforms. To fully capture the changes in live content, three challenges must be addressed: (1) As depicted in Figure 1 (a), live-streaming scenarios involve several modalities, including real-time visual frames, speech from the streamer, bullet comments from audiences, and other categorical inputs (e.g., streamer IDs and live IDs). Consequently, a model that can leverage various modalities and their historical information is crucial. (2) Traditional pointwise estimation may lead to the model converging on a sub-optimal solution. Apart from that, as illustrated in Figure 1 (b), the traditional ID based model may learn a mediocre solution (shown in green line), neglecting multiple singular points that could represent live streaming highlights. Thus, an appropriate loss function that models the contrastive relationship between highlight and non-highlight frames is beneficial for generating the final CTR prediction results. (3) In live-streaming scenarios, the streamer’s speech and audiences’ comments may be ambiguous and not sequentially aligned with the visual frames, necessitating a module to filter out the noise caused by misalignment.

In this paper, we propose an end-to-end transformer-based network called ContentCTR, which leverages multi-modality features for frame-level Click-Through Rate (CTR) prediction. To the best of our knowledge, this is the first research that explores frame-level live streaming recommendation. First, the ContentCTR exploits the

information contained in visual, speech, comment, and ID embedding and uses an attention mechanism to adaptively find correlations and interactions from different modalities. Second, to alleviate the misalignment between visual and textual modality, we develop a dynamic time warping algorithm to address potential temporal discrepancies that may arise during live streaming events. Finally, to better exploit the contrastive information between highlight frames and no-highlight frames, we design a pairwise loss function with first-order difference constraints. We find that the constraints are essential when jointly optimizing pointwise and pairwise losses to avoid collisions and model collapse. It is worth noting that our work focuses on modeling the relationship between real-time live content and CTR, without introducing user personalization factors. The model can be utilized not only for online traffic mechanism but also for helping existing recommendation models enhance content understanding. We perform comprehensive experiments on both a large-scale real-world live streaming dataset and a public PHD dataset [10] and achieve the state-of-the-art performance.

In summary, the main contributions made in this work are as follows:

- We propose ContentCTR, an end-to-end transformer-based network for frame-level CTR prediction in live streaming platform. It efficiently utilizes the features of different modalities and captures the dynamic highlight pattern.
- We design a pairwise-based loss function with first-order difference constraints to exploit the contrastive information of highlight frames and no-highlight frames. In addition, a dynamic time warping (DTW) based alignment strategy is present to alleviate misalignment in streaming scenarios.
- We conduct comprehensive experiments on both a large-scale real-world live streaming dataset and a public PHD dataset [10]. Our method outperforms all baseline models. We also perform online A/B testing on the live streaming platform of company, achieving a 2.9% lift in CTR and a 5.9% improvement in terms of live play duration.

2 RELATED WORK

2.1 Deep Click-Through Rate Prediction

Based on the input modality, the existing methods for Click-Through Rate (CTR) prediction can be classified into two categories: ID-based and multi-modal-based CTR prediction models. For the ID-based CTR prediction approaches, techniques such as Wide&Deep [6], DeepFM [11], and DCN-M [31] aim to capture high-order feature interactions and complex feature dependencies. On the other hand, DIFM [19] addresses the fixed feature representation problem by learning the input vector level weights for feature representations. For multi-modal-based CTR prediction models, the majority of the research has focused on video Click-Through Rates prediction, which is a fundamental task in the field of multimedia and information retrieval. For example, AutoFIS [16] propose a two-stage algorithm called Automatic Feature Interaction Selection while UBR4CTR [24] propose User Behavior Retrieval for CTR prediction framework to tackle the problem that sequential patterns such as periodicity or long-term dependency are not embedded in the recent several behaviors but in far back history. Both of these approaches are based on the Factorization Machine (FM) [26] method.

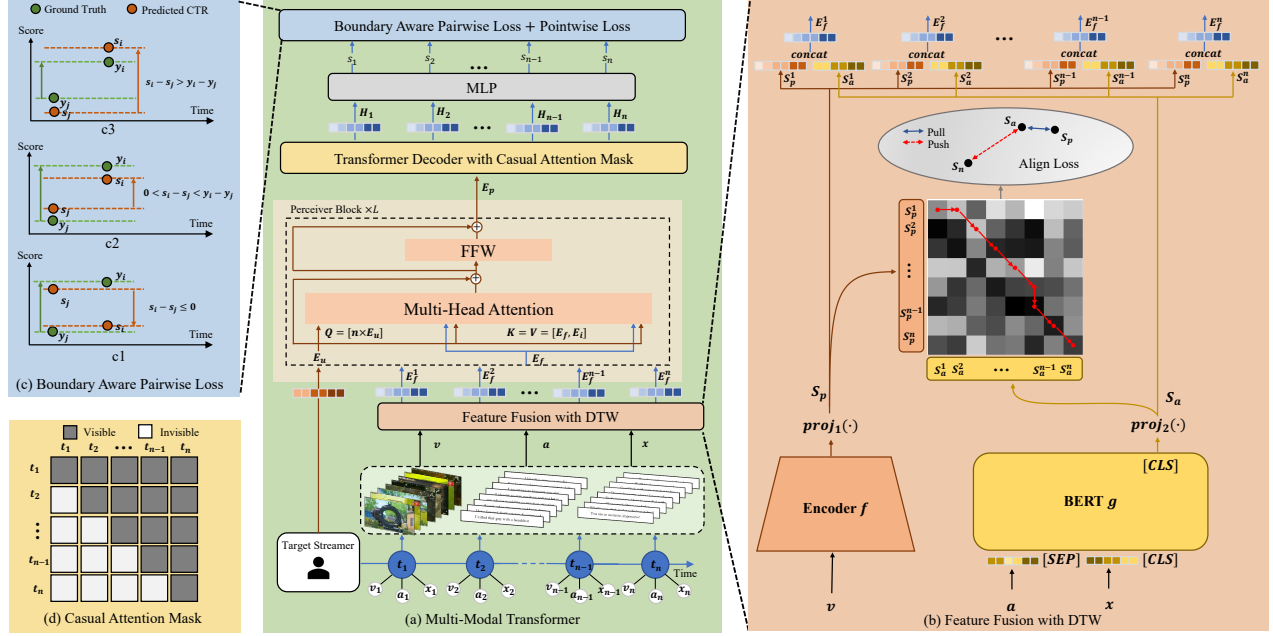


Figure 2: The framework of our proposed method. Part (a) and (d) show the architecture of Multi-Modal Transformer and Casual Attention Mask which are discussed in Section 3.2. Part (b) shows the proposed Dynamic Time Warping Alignment Loss which is discussed in Section 3.3. Part (c) shows the motivation of Boundary-aware Pairwise Loss which is discussed in Section 3.4.

For DNNs based methods, DSTN [22] propose a Deep Spatiotemporal Translation Network (DSTN), which utilize the auxiliary data in spatial and temporal domain. In order to better capture user preferences, HyperCTR [12] design a Hypergraph Click-Through Rate prediction framework, built upon the hyperedge notion of hypergraph neural networks, which can yield modal-specific representations of users and micro-videos. However, all these methods only focus on video-level CTR prediction. Compared with traditional video CTR prediction task, stream highlight CTR prediction needs frames level CTR prediction network which can leverage multi-modal feature to address its unique challenges, i.e. the dynamic patterns of live-streaming and non-sequential alignment of different modality.

2.2 Live Streaming Recommendation

As an emerging form of social media, increasing attention has being given to live streaming recommendations. Current works typically treat streamers or audiences as items, and develop recommendation systems to predict interactions between users and items. These methods can be broadly categorized into two types: Filtering-based methods and Deep learning-based methods. For Filtering-based methods, HyPAR [33] proposed a Collaborative and Content-based filtering hybrid approach to recommend streamers to audiences. For Deep learning-based methods, LiveRec [25] and [35] leverages self-attentive mechanisms based on historical interactions and repeat consumption behavior, while DRIVER [9] learns dynamic representations by leveraging users' highly dynamic behaviors. [37] utilizes LSTM to extract the representation and capture the preference of streamers and audiences. However, these works only focus on the

interaction between streamers and audiences at the item level, neglecting the importance of highly dynamic multi-modal content in live streaming. Additionally, they cannot provide frame-level recommendations.

2.3 Video Highlight Detection

The task most closely related to stream CTR prediction is Personalized Video Highlight Detection (P-VHD), as both aim to recognize different content patterns in temporal sequences. The primary objective of P-VHD is retrieving a subset of video frames that capture a person's main attention or interests from the original video. Recently, several works [2, 5, 10, 27, 29] are proposed for P-VHD task, which aims at extracting user-adaptive highlight predictions guided by annotated user history. For example, PHD-GIFs [10] is the first personalized video highlight detection technique that also creates a large-scale dataset called PHD. The P-VHD task and stream CTR prediction both need to consider personalized behavior, while P-VHD only needs to predict the label of frames to highlight or no-highlight one. On the other hand, stream CTR prediction needs to regress the correct order of all frames, which means that stream highlight CTR prediction is more challenging.

3 METHODOLOGY

As illustrated in Figure 2, we propose the Transformer Click-Through Rate prediction framework, named ContentCTR. ContentCTR utilizes multi-modal features to predict the CTR at frame level. In particular, we first define the problem of predicting CTR for stream highlights in Section 3.1. Then the Multi-Modal Transformer backbone is introduced to fuse and interact with different modalities in

Section 3.2. Additionally, we propose the Dynamic Time Warping strategy for aligning text and visual feature in Section 3.3. Finally, boundary-aware pairwise loss for exploring contrasting information is presented in Section 3.4.

3.1 Problem Formulation

The main goal of the CTR prediction is to estimate the probability of the current content of the live broadcast being clicked. In this work, we simplify this problem for predicting the Click-Through Rate of frame δ_i at timestamp i . We denote $M_i = \{v_i, a_i, x_i\}$ as the multi-modal tuple, where v_i, a_i and x_i represent the visual, speech and comments of frame δ_i , respectively. We hypothesize that different streamers have distinct talents and attract different audiences who are typically interested in specific types of highlight moments in a streaming room. For instance, as shown in Figure 1 (b), some audiences may be attracted to a streamer's dancing, while others may be attracted to the PK between streamers. Therefore, we denote E_u as the ID embedding of streamer u . The ID embedding of streamer is extracted by a pre-trained SIM [23] model that reflects the rough classification of a streamer. Additionally, we consider that the highlight pattern and audience taste change over time, and thus, the model should use information from the $n - 1$ lookahead windows $W_M = \{M_{i-n+1}, M_{i-n+2}, \dots, M_{i-1}\}$ of previous frames to predict the CTR y_i of frame δ_i . The prediction problem can be formulated as follows,

$$\text{Prob}(\delta_i | W_M, E_u) \sim \Gamma(W_M, E_u, \delta_i) \quad (1)$$

where the frame δ_i is represented by multi-modal feature of lookahead windows W_M and the ID embedding E_u . The probability that audiences will click on the frame δ_i is denoted by $\text{Prob}(\delta_i | W_M, E_u)$. Moreover, $\Gamma(W_M, E_u, \delta_i)$ is the model used to estimate the probability $\text{Prob}(\delta_i | W_M, E_u)$.

3.2 Multi-Modal Transformer

The proposed Multi-Modal Transformer backbone includes the following three basic components, i.e., feature fusion layer, Perceiver block, and casual sequence decoder.

3.2.1 Feature Fusion Layer. As depicted in Figure 2 (b), given the historical window W_M of streamer u , we extract multi-modal features for every timestamp, including the streaming frames v_i , Auto Speech Recognition (ASR) text a_i from the streamer, comment x_i from the audiences. The streaming frames are tokenized by the pre-trained swin [18] $f(\cdot)$ while the ASR and comment text are tokenized with the BERT [7] Chinese Large model $g(\cdot)$ and we only use the hidden feature of <CLS> token as the text embedding. Since the embedding dimensions for streaming frames and language tokens are different, two MLP heads $proj_1(\cdot)$ and $proj_2(\cdot)$ are set to map the frames embedding and text embedding to the same dimension d , denoted by $S_p \in \mathcal{R}^{b \times n \times 1 \times d}$ and $S_a \in \mathcal{R}^{b \times n \times 1 \times d}$, where n is the length of lookahead window and b is the batch size. The above process can be formulated as follows:

$$\begin{aligned} S_p &= proj_1(f(v)) \\ S_a &= proj_2(g(\text{CLS})([a, x])) \end{aligned} \quad (2)$$

where $v = [v_1, \dots, v_n]$, $a = [a_1, \dots, a_n]$ and $x = [x_1, \dots, x_n]$.

For frames δ_i at timestamp i , the corresponding frames feature and text feature are denoted by $S_p^i \in \mathcal{R}^{b \times 1 \times 1 \times d}$ and $S_a^i \in \mathcal{R}^{b \times 1 \times 1 \times d}$ respectively. We concatenate S_p^i and S_a^i as the final input tokens $E_f^i \in \mathcal{R}^{b \times 1 \times 2 \times d}$ of the Perceiver block. Finally, we obtain $E_f \in \mathcal{R}^{b \times n \times 2 \times d}$ by concatenating all features of the timestamps.

3.2.2 Perceiver Block. The Perceiver block is connected to the casual sequence decoder, as shown in Figure 2 (a). It takes the repeated ID embedding $n \times E_u \in \mathcal{R}^{b \times n \times 1 \times d}$ of streamer u and the fusion multi-modal feature $E_f \in \mathcal{R}^{b \times n \times 2 \times d}$ as input, where n is the length of lookahead window, b is the batch size and d is the input dimension. The motivation of Perceiver block is the success of multi-head attention mechanism [28] in sequential recommendation [4] and video understanding [1]. We apply it to capture streamer's highlight patterns on lookahead sequences with the query-based retrieval problem.

First, we initialize learned latent query $Q \in \mathcal{R}^{bn \times 1 \times d}$ with flattened $n \times E_u$. Next, we concatenate the flattened E_f and $n \times E_u$ at the second dimension and take $K = V = [E_f, n \times E_u] \in \mathcal{R}^{bn \times 3 \times d}$ as key and value. We linearly project the input vector to latent vectors with d_h dimensions, by different linear projections. Then we perform the scaled dot-product multi-head attention as follows,

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_{n_h}), \\ \text{head}_i &= \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right), \end{aligned} \quad (3)$$

where n_h is the number of attention head and $W_i^Q \in \mathcal{R}^{d \times d_h}$, $W_i^K \in \mathcal{R}^{d \times d_h}$ and $W_i^V \in \mathcal{R}^{d \times d_h}$ are the learnable parameters. The scaled dot-product attention function is defined as follows,

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right) \quad (4)$$

Then, the Perceiver block employs a feed-forward network with residual connections for performance boosting. The Perceiver block is stacked for L layers and the outputs of previous block are fed into the Perceiver block alternately. The output of Perceiver block is denoted by $E_p \in \mathcal{R}^{bn \times 1 \times d_n}$. The pseudocode of the Perceiver block is shown in Algorithm 1.

Algorithm 1 A Pytorch-style Pseudocode for Perceiver Block.

```
def perceiver_block(
    x_f, # The [b, n, 2, d] multi-modal feature
    x, # The learned latent query with shape [b, n, 1, d]
    num_layers, # The number of layers
):
    x_f = flatten(x_f) # [b, n, 2, d] -> [b * n, 2, d]
    x = flatten(x) # [b, n, 1, d] -> [b * n, 1, d]
    for i in range(num_layers):
        # Attention
        x = x + attention_i(q=x, kv=concat([x_f, x]))
        # Feed forward with residual connection
        x = x + ffw_i(x)
    return x
```

3.2.3 Casual Sequence Decoder. We first unflatten and squeeze the output of Perceiver block $E_p \in \mathbb{R}^{bn \times 1 \times d}$ as the dimension of $\mathbb{R}^{b \times n \times d}$. Then we initialize the query, key and value of decoder as $Q = K = V = E_p \in \mathbb{R}^{b \times n \times d}$. After that, we apply the scaled dot-product multi-head attention which is the same as shown in Equation 3, but with the scaled dot-product attention function and casual attention as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}} + M\right) \quad (5)$$

where M is the casual attention mask. As shown in Figure 2 (d), it is an $n \times n$ matrix filled with -inf and its upper triangular sub-matrices are filled with 0. By applying the casual attention mask M the problem of future information leakage in the temporal dimension is avoided. The pseudocode of the Decoder is shown in Algorithm 2. The output H from the final attention layer is then fed into a

Algorithm 2 A Pytorch-style Pseudocode for Decoder.

```
def decoder(
    x, # The input with shape [b*n, 1, d_n]
    num_layers, # The number of layers
):
    x = unflatten(x).squeeze() # [b*n, 1, d_n] -> [b, n, d_n]
    for i in range(num_layers):
        # Attention
        x = x + attention_with_mask_i(q=x, k=x, v=x)
        # Feed forward with residual connection
        x = x + ffw_i(x)
    return x
```

fully-connected layer, followed by a sigmoid transformation to produce the scalar prediction of CTR s . The parameters of the fully-connected layer are represented by $W \in \mathbb{R}^{d_n \times 1}$.

The main objective is to maximize the log-likelihood between the predicted CTR s and the actual CTR y , which is achieved using the following pointwise model to optimize the standard LogLoss [17]:

$$L_{\text{Point}} = -\frac{1}{n} \sum_{i=1}^n [y_i \cdot \log(s_i) + (1 - y_i) \cdot \log(1 - s_i)] \quad (6)$$

3.3 DTW Alignment

The motivation behind for alignment is to address potential temporal discrepancies that may arise during live streaming events. For example, the streamer may describe the plan before taking action, or explain detailed information after action. Additionally, the comments from the audiences may experience some time lag, which further exacerbates the misalignment issue. Therefore, it is essential to train text and visual encoders that can handle misalignment to alleviate that problem. Inspired by previous works [15, 32], in this section we present our contrastive learning based framework for visual and text sequences alignment.

Given the text sequence S_a and visual sequence S_p from a look-ahead window W_M , we assume that when timestamps i in the text sequence and j in the visual sequence represent the same pattern, then S_a^i and S_p^j should have a semantic similarity. To train a network that

can minimized the distance between text sequences S_a and visual sequences S_p , we use the Dynamic Time Warping [20] (DTW) which calculates the minimum cumulative matching costs over units as the sequence distance to measure the similarity. First, we computes a pairwise distance matrix $D(S_a, S_p) := \left[\text{sim}(S_a^i, S_p^j) \right]_{ij} \in \mathbb{R}^{n \times n}$ with a distance measure $\text{sim}(\cdot)$. In our work we applied the cosine similarity as $\text{sim}(\cdot)$. Then, we employs dynamic programming and sets a matrix $C \in \mathbb{R}^{n \times n}$ to record the minimum cumulative cost between S_a^i and S_p^j [8]:

$$C_{i,j} = D_{i,j} + \min \{C_{i-1,j-1}, C_{i-1,j}, C_{i,j-1}\} \quad (7)$$

where $1 \leq i, j \leq n$. Then, the distance $d_{\{S_a, S_p\}}$ between sequences S_a and S_p is set to the last element of matrix C :

$$d_{\{S_a, S_p\}} = C_{n,n} \quad (8)$$

Contrastive Learning (CL) is a self-supervised learning technique that learns a representation of data by comparing different views of same samples. The basic idea of CL is to bring together similar pairs and push away dissimilar pairs. We hypothesize that the positive pair $\{S_a, S_p\}$ should be similar. To construct the negative sample S_n , we randomly shuffle the temporal order of video sequence S_p . Therefore, the negative pair $\{S_a, S_n\}$ should be dissimilar. Then, we can derive the training objective to minimize the InfoNCE loss [21]:

$$\mathcal{L}_{\text{align}}(S_a, S_p, S_n) = -\log \frac{\exp(d_{\{S_a, S_p\}}/\tau)}{\exp(d_{\{S_a, S_p\}}/\tau) + \sum_{S_n \in S_n} \exp(d_{\{S_a, S_n\}}/\tau)} \quad (9)$$

where $S_n = \{S_{n_i}\}_{i=1}^N$ is the set of N negative samples shuffled from S_p . As shown in Figure 2 (b), by minimizing the align loss $\mathcal{L}_{\text{align}}$, the visual encoder $f(\cdot)$ and text encoder $g(\cdot)$ should be able to learn a good representation from aligned and misaligned pairs.

3.4 Boundary Aware Pairwise Loss

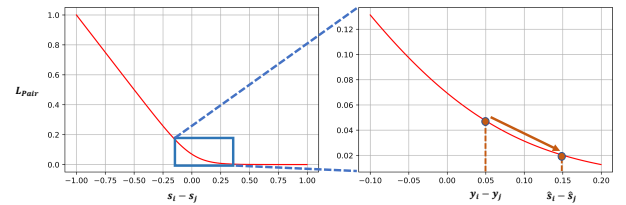


Figure 3: Loss function for pairwise optimization without any constraints.

In this section, we present some intriguing discoveries on the traditional pairwise logistic ranking loss [3] and propose modifications to the loss function by integrating the boundary-aware first-order difference constraint, thereby enhancing the optimization. The addition of the pairwise loss aids in exploiting the underlying contrasting information present in both the highlight and no-highlight frames.

Consider the following pairwise loss function, which has no constraints:

$$L_{\text{Pair}}^0 = \sum_{y_i > y_j} \log \left(1 + e^{-\sigma(s_i - s_j)} \right) \quad (10)$$

where y_i and y_j represent the ground truth CTR at timestamps i and j , while s_i and s_j represent the predicted CTR from the model.

Figure 3 illustrates the changing of the loss function L_{Pair} with respect to $s_i - s_j$, which reveals that when $y_i - y_j$ holds, minimizing L_{Pair} tends to cause $s_i - s_j$ to overtake the optimal value of $y_i - y_j$. This leads to over-optimization.

Based on above findings, we propose a revised pairwise loss that incorporates a boundary-aware first-order difference constraint:

$$L_{Pair}^1 = \sum_{y_i > y_j} \log(1 + e^{-\sigma(s_i - s_j)}), (y_i - y_j) - (s_i - s_j) \geq 0 \quad (11)$$

where $(y_i - y_j) - (s_i - s_j) \geq 0$ denotes the boundary, and solely those samples that reside within the boundary will calculate L_{Pair}^1 . Any samples located outside the boundary will result in L_{Pair}^1 being set to 0.

Without losing generality, the original pairwise loss function L_{Pair}^0 presented in Equation 10 can be divided into three distinct parts:

- Part 1: As shown in Figure 2 (c1), when $s_i - s_j \leq 0$, the model's assessment of the significance between timestamp i and j is entirely incorrect, given that timestamp i is more "highlighting" than timestamp j .
- Part 2: As shown in Figure 2 (c2), when $0 < s_i - s_j < y_i - y_j$, it implies that the model has distinguished that timestamp i is more "highlighting" than timestamp j , but it still fails to accurately predict the difference in CTR values between the two timestamps, and thus, it is still not optimal.
- Part 3: As shown in Figure 2 (c3), when $y_i - y_j < s_i - s_j$, it indicates that the model's predicted CTR value is too aggressive, resulting in over-optimization.

In order to verify above scenarios, we design the following loss functions:

$$L_{Pair}^2 = \sum_{y_i > y_j} \log(1 + e^{-\sigma(s_i - s_j)}), s_i - s_j \leq 0 \quad (12)$$

where $s_i - s_j \leq 0$ means that it only optimizes on Part 1.

$$L_{Pair}^3 = \sum_{y_i > y_j} \log(1 + e^{-\sigma(s_i - s_j)}), y_i - y_j > s_i - s_j > 0 \quad (13)$$

where $y_i - y_j > s_i - s_j > 0$ means that it only optimizes on Part 2. The ablation study among L_{Pair}^0 , L_{Pair}^1 , L_{Pair}^2 and L_{Pair}^3 is discussed in Section 4. In this work, we apply the L_{Pair}^1 for optimization.

By combining pointwise loss, align loss and pairwise loss, our final loss used to learn the model parameters is defined as:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{Point} + \lambda_2 \mathcal{L}_{align} + \lambda_3 \mathcal{L}_{Pair}^1 \quad (14)$$

where λ_1 , λ_2 and λ_3 are the tradeoff parameters.

4 EXPERIMENTS

In this section, we first introduce the datasets, evaluation metrics, and implementation details. Next, we present the experimental results and some analysis of them. Specifically, our experiments aim to answer the following research questions:

- **RQ1:** How does the modality impact the model?
- **RQ2:** To what extent can the proposed Perceiver Block enhance the model's performance?
- **RQ3:** How effective is the proposed DTW alignment block in alleviating misalignment cases?

- **RQ4:** How much can the model's performance be enhanced by the proposed pairwise loss?
- **RQ5:** How well does the proposed model perform on public datasets, such as video highlights?

4.1 Dataset

To comprehensively evaluate the model's performance, we present experimental results on both KLive Dataset and a public video highlighting dataset. Detailed information regarding these datasets is provided as follows.

4.1.1 KLive Dataset. We have constructed a large-scale dataset based on our company's short video and live streaming platform, which boasts over three million daily active users. Our methodology involved selecting 39,000 high-quality live rooms over a period of three days, based on user rewards and viewing time indicators. Each live room is then divided into multiple consecutive 30s live segments, with three pictures evenly sampled for each segment. The streamer's ASR and audiences' comments are extracted for each segment, and the CTR is calculated by dividing the number of clicked users by all watched users. We then construct consecutive 20 streaming segments into a dataset sample, filtering out samples with a low number of watched users of the last live segment to ensure a reliable CTR. In the end, we obtained 1,436,979 and 286,510 samples for our training and test datasets, respectively. Each sample is comprised of the streamer's item ID, three video frames for 20 streaming segments, audience comments, the streamer's ASR speech, and the ground truth CTR for the all live segments.

4.1.2 PHD Dataset. To verify the generality of our method, we evaluate on the publicly available personalized video highlight detection dataset [10] (PHD). This dataset comprises of URLs of YouTube videos, IDs of evaluators or "users," and the segments they designated as highlight frames based on their preferences. The most recent video that a user annotated is considered as the target video for that particular user. Since PHD only provides YouTube URLs, we download the original videos and crawl the captions from YouTube to carry out the experiments. Our training and testing sets include up to 20 history highlight frames and one target video per user. The duration of the history highlight frames varies from 1 to 672.19 seconds, with an average length of 5.12 seconds. The target videos range from 1 to 37,434 seconds, with an average length of 431.79 seconds. The training set comprises of a total of 12,541 users and 81,056 videos, while the testing set has 833 users and 7,595 videos that do not overlap with any user or video in the training set. Consistent with [5], we divide the target video into fixed-length segments (192 frames) and only train those segments that contain highlight frames. During testing, we make the entire video segments as input for inference.

4.2 Evaluation Metrics

For the experiments on KLive dataset, we employed Kendall's tau τ [14] to measure the correlation between our predicted click-through rates (CTR) s and the ground truth CTR y . It is defined as:

$$\tau = \frac{P - Q}{\sqrt{(P + Q + T) \cdot (P + Q + U)}} \quad (15)$$

where P represents the number of concordant pairs, Q denotes the number of discordant pairs, T indicates the number of ties only in s , and U is the number of ties only in y . If the same pair experiences a tie in both s and y , it is not included in either T or U . Value of τ close to 1 indicate strong agreement, while values approaching -1 indicate strong disagreement.

Regarding the experiment on the PHD dataset, we utilized the widely adopted mean Average Precision (mAP) as a metric to evaluate the performance of our method, which is also applied in previous works [2, 5, 27, 29] in video highlight detection. We report the mAP on the test set and follow the way in [29] to calculate the mAP.

4.3 Implementation Details

4.3.1 Baseline: On KLive dataset, we compare our method with ID-based recommendation methods and several variants of our proposed approach as follows:

- **W&D [6]:** This method trains wide linear models and deep learning network together to integrate the advantages of memorization and generalization in recommendation, but requires feature engineering in addition to raw features.
- **DeepFM [11]:** This method utilizes the benefits of both factorization machines and deep learning by sharing input for the deep and wide components.
- **DIFM [19]:** This approach leverages the benefit of input-aware factorization machines to address the issue of standard FMs that generate a fixed representation for varying input instances.
- **DCN-M [31]:** This approach is an enhanced version of DCN [30], with greater effectiveness in learning both explicit and implicit feature crosses through the application of low-rank techniques to approximate such crosses.
- **SelfAttention-plain:** We remove the ID embedding E_u of the streamer of ContentCTR and replace the Perceiver block with a self-attention block.
- **SelfAttention-input:** We replace the Perceiver block with a self-attention block and concatenate the ID embedding E_u of the streamer with E_f to be used as input for the self-attention.
- **CrossAttention-q:** We replace the Perceiver block with a cross-attention block. The query of every cross-attention layer is set as E_u , while the key and value are initialized with E_f .
- **CrossAttention-kv:** We replace the Perceiver block with a cross-attention block. The key and value of every cross-attention layer are set as E_u , while the query is initialized with E_f .
- **CrossAttention-qkv:** We replace the Perceiver block with a cross-attention block. The query is initialized with E_u , while the key and value are initialized with E_f .

On the PHD dataset, we make slight modifications to our proposed ContentCTR method and compare it with several classic baselines for personalized video highlight detection. Specifically, they are:

- **Adaptive-H-FCSN [27]:** This method employs a convolutional highlight detection network with a history encoder to learn user-specific highlight patterns.
- **PR-Net [5]:** This method proposes a reasoning framework to explicitly learn frame-level patterns. It also employs contrastive learning to alleviate annotation ambiguity.

- **PAC-Net [29]:** This method introduces a Decision Boundary Customizer (DBC) module and a Mini-History (Mi-Hi) mechanism to capture more fine-grained user-specific preferences.
- **ShowMe [2]:** This method leverages the content of both user history and target videos, using pre-trained features of YOLOv5 [13].
- **ContentCTR-HL:** We modify our proposed ContentCTR model to adapt it to the personalized video highlight detection task. First, we extract YOLOv5 features from images for training and testing. The text comes from the captions crawled from the YouTube video. Second, the ID embedding is replaced with the user history feature for personalized highlighting. Then, since there is no need to consider future information leakage in the video highlighting task, the causal attention mask is set to fully visible. Finally, when optimizing the network, only \mathcal{L}_{Point} and \mathcal{L}_{align} are reserved for optimization.

4.3.2 Setup Detail: During the training of ContentCTR on the KLive dataset, the layer number for the Perceiver block and Decoder block is set to 3, the input dimension $d = 512$, the hidden dimension of the transformer d_h is set to 64 and the number of attention head $n_h = 8$. We utilize pre-trained swin as the Encoder $f(\cdot)$ and pre-trained Bert Chinese Large as the Encoder $g(\cdot)$. The swin and BERT are both pre-trained on image and text data from 39,000 streaming rooms. The tradeoff parameters, λ_1 , λ_2 , and λ_3 , are respectively set to 0.65, 0.15, and 0.20. The σ in L_{Pair}^1 is set to 10 and the number of negative sample N is set to 8. We optimize ContentCTR for 12 epochs with a learning rate of 5×10^{-5} using the Adam optimizer, and the global batch size is set to 48. During training, the gradient of swin and BERT feature extractors are open for fine-tuning, while the parameters of two MLP heads are updated, too. The training process takes about 34 hours on 8 Nvidia Tesla V100 GPUs. During the training and testing of ID-based recommendation methods on the KLive dataset, the input sparse features for these methods include the streamer item ID, live ID, timestamp ID, exposure count, comment count, gift count, click count, like count, follow count, room entry and exit count for each streaming segment. The objective of these methods is to regress the CTR of each segment. When training ContentCTR-HL on the PHD dataset, we first extract the YOLOv5 features for each frame of the historical highlight frames in the training and testing set. We use pre-trained YOLOv5 as the Encoder $f(\cdot)$ and pre-trained Bert English Large as the Encoder $g(\cdot)$. Then we train our network using the Adam optimizer with a batch size of 8 and an initial learning rate of 5×10^{-5} , with a cosine decay scheduler. The training is performed for 20 epochs on 8 NVIDIA Tesla V100 GPUs, which takes approximately 16 hours.

4.4 Overall Performance Comparison

4.4.1 Quantitative Results: Table 1 summarizes the CTR prediction performances achieved by various methods on KLive dataset, indicating that **ContentCTR** surpasses all ID-based recommendation methods and the variations of **ContentCTR**. The findings reveal that the interaction approach among different modalities plays a crucial role in the final CTR prediction performance. Furthermore, compared to ID-based recommendation methods, ContentCTR demonstrates superior CTR prediction performance which

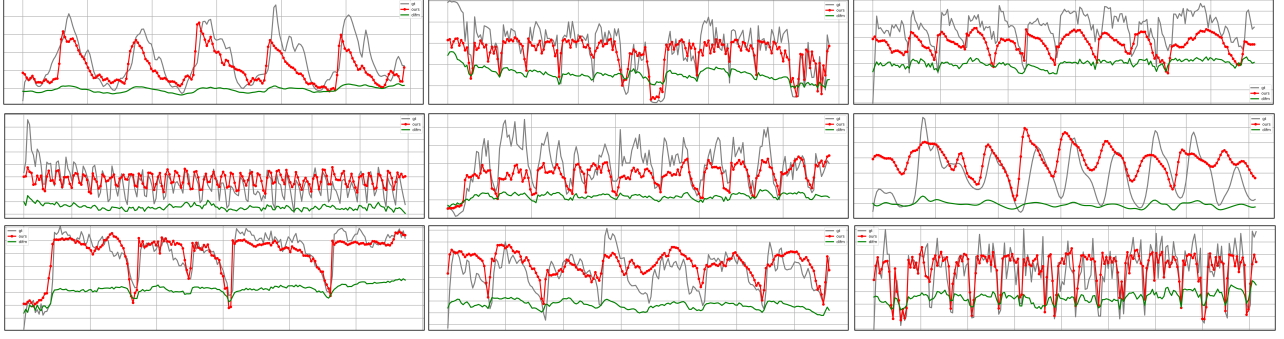


Figure 4: Visualization of the CTR prediction results for ContentCTR (red line), DIFM (green line) and ground truth CTR (grey line).

Table 1: Performances of different methods on KLive dataset

Model		Tau τ
<i>ID-based Recommendation Methods</i>		
W&D	[DLRS'16]	0.5619
DeepFM	[IJCAI'17]	0.5542
DIFM	[IJCAI'21]	0.5697
DCN-M	[WWW'21]	0.5629
<i>Variants of ContentCTR</i>		
SelfAttention-plain		0.5868
SelfAttention-input		0.5718
CrossAttention-q		0.5909
CrossAttention-kv		0.5780
CrossAttention-qkv		0.5883
ContentCTR		0.5919

suggests that the transformer based encoder-decoder architecture and Perceiver block is more efficient in capturing dynamic patterns.

Table 2: Comparison with state-of-the-art alternative on PHD dataset.

Model		Param.(M)	mAP(%)
Adaptive-H-FCSN [27]	[ECCV'20]	197.35	15.65
PR-Net [5]	[ICCV'21]	-	18.66
PAC-Net [29]	[ECCV'22]	5.89	17.51
ShowMe [2]	[MM'22]	-	16.40
ContentCTR-HL	[Ours]	66.35	21.89

In order to answer **RQ5** and further evaluate the effectiveness of our method, we report the performance comparison of ContentCTR-HL with various state-of-the-art alternatives on PHD dataset. As shown in Table 2, our approach outperforms all other alternatives on PHD dataset, outperforming **PR-Net** [5] by +3.23%. This result undoubtedly demonstrates the practicality and effectiveness of our method. The key improvement of our method can be derived from the following three aspects. Firstly, in contrast to other approaches that solely rely on visual features, our method leverages the multi-modal features of video frames and captions. Secondly,

the sequence-to-sequence encoder-decoder architecture is proficient in detecting frame-level highlights. Thirdly, the incorporation of DTW alignment effectively alleviates the possible misalignments between captions and video frames.

4.4.2 Qualitative Results: Moreover, we visualize several typical cases on KLive dataset in Figure 4. Figure 4 illustrates the prediction results of **ContentCTR** in red line, **DIFM** in green line and the ground truth CTR in grey line. It is obvious that when the ground truth CTR arises abrupt fluctuations (shark rises or fails) due to the dynamic changing of content, **ContentCTR** is able to capture real-time content changes more effectively. However, **DIFM** tends to change more gradually, which is improper for high dynamic streaming highlight recommendation scenarios.

4.5 Modality Impact(RQ1)

We study the impact of different modalities on ContentCTR's performance on KLive dataset, as shown in Table 3. The results demonstrate that ContentCTR achieves a tau of 0.5919 when all the modalities are engaged. We find that the visual modality has the most important impact on the model performance, causing a performance degradation of -4.72% when removed. Apart from that, the text modality is the second most significant factor, leading to a -2.01% degradation. Lastly, the ID embedding has the smallest but still significant effect on the model's performance, with a -0.51% degradation when removed. This ablation study suggests that incorporating different modalities is essential for the model to predict CTR accurately.

Table 3: Ablation study on different modality impact. We study the effect of visual modality v , speech modality a , comment modality x and item modality u .

Model	v	a	x	u	Tau τ
ContentCTR	✓	✓	✓	✓	0.5919
ContentCTR w/o item	✓	✓	✓	-	0.5868 ↓ 0.51%
ContentCTR w/o text	✓	-	-	✓	0.5718 ↓ 2.01%
ContentCTR w/o visual	-	✓	✓	✓	0.5447 ↓ 4.72%

4.6 Perceiver Block(RQ2)

As demonstrated in Table 1, the Perceiver block surpasses all the variants of ContentCTR in terms of performance. Table 1 provides the following insights.

Firstly, the results indicate that inputting the ID embedding directly into the transformer model fails to fully leverage personalized information. The self-attention architecture is not suitable for detecting the correlation between different modalities.

Secondly, concerning cross-attention architectures, it is critical to carefully design the interaction approach between the ID embedding and other content-based modalities. The traditional query-key-value-based cross-attention mechanism proves to be less effective than the Perceiver block.

4.7 DTW Alignment and Pairwise Loss(RQ3&RQ4)

We design **Model1-Model5** to investigate the impact of different loss functions on KLive dataset, which include L_{Point} , L_{Pair}^0 , L_{Pair}^1 , L_{Pair}^2 , L_{Pair}^3 , and L_{align} . According to Table 4, when we optimized **Model4** with L_{Point} and L_{Pair}^2 , its performance significantly degraded by -5.05%. We hypothesize that the gradient changes are too drastic when optimizing only Part 1, which is detrimental to modeling the contrastive information of highlight frames and non-highlight frames. However, when optimizing only Part 2 with L_{Pair}^3 , **Model5** showed a performance improvement of +0.66%, but it still falls behind **Model3**, which is jointly optimized on Part 1 and Part 2 with L_{Pair}^0 , resulting in a significant improvement of +1.11%. Although **Model2** with L_{Pair}^0 has also shown performance

Table 4: Ablation study of ContentCTR on pairwise loss and DTW align loss. We study the effect of loss function L_{Point} , L_{Pair}^0 , L_{Pair}^1 , L_{Pair}^2 , L_{Pair}^3 and L_{align} .

Model	L_{Point}	L_{Pair}^0	L_{Pair}^1	L_{Pair}^2	L_{Pair}^3	L_{align}	Tau τ	avg. s/y
Model1	✓	-	-	-	-	-	0.5761	1.2789
Model2	✓	✓	-	-	-	-	0.5857 \uparrow 0.96%	1.3164
Model3	✓	-	✓	-	-	-	0.5872 \uparrow 1.11%	1.3021
Model4	✓	-	-	✓	-	-	0.5256 \downarrow 5.05%	1.3414
Model5	✓	-	-	-	✓	-	0.5824 \uparrow 0.66%	1.2194
Ours	✓	-	✓	-	-	✓	0.5919 \uparrow 1.58%	1.3011

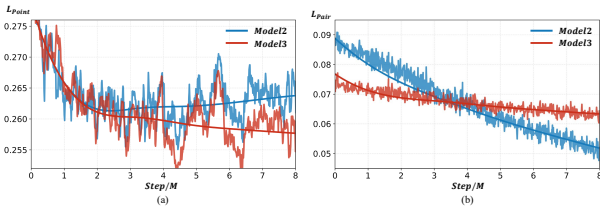


Figure 5: The change of L_{Point} and L_{Pair} of Model2 (blue line) and Model3 (red line) during training. It is obvious that the pairwise loss of Model2 is over-optimized which causes the pointwise loss to collapse in Model2 while both losses remain normal in Model3.

improvement, we have noticed that during training, the pointwise

loss tends to collapse due to the over-optimization of pairwise loss on Part 3, which is shown in Figure 5.

As expected, when **Model3** incorporates L_{align} for additional optimization, the DTW alignment loss shows further improvement by +0.47%. Figure 6 shows some visualizations of DTW alignment results for the KLive dataset, which reveal that misalignment cases do exist in streaming scenarios. This also indicates that the involvement of L_{align} helps ContentCTR in training better visual and text encoders which reduces the possible misalignment between the two.

4.8 Ablation Study on PHD

We also conducted an additional ablation study on the PHD dataset based on ContentCTR-HL, analyzing the impact of highlight frames, captions, DTW align loss and input frames number. As shown in Table 5, the performance of ContentCTR-HL decreased by -2.34% when the history feature is removed, indicating that the history feature of highlight frames is crucial for personalized video highlight detection. Similarly, the absence of captions resulted in a -1.11% decrease in ContentCTR-HL's performance, highlighting the usefulness of text modality for the model's detection capability. Interestingly, ContentCTR-HL is found to be less sensitive to DTW alignment loss. We hypothesize that this phenomenon is caused because most captions generated by the official YouTuber are accurate enough. Apart from that, since L_{align} is a self-learning contrastive loss, it may require a larger data scale to take effect, but currently, only 30% of all videos in PHD have downloadable captions. We also find that the performance of ContentCTR-HL improves with the increase of the input frames number. However, the increase in performance is insignificant when the number of frames is increased from 96 to 192.

Table 5: Ablation Study on ContentCTR-HL.

Model	Frames Number n			
	32	64	96	192
ContentCTR-HL(Ours)	18.29	19.54	21.59	21.89
Ours w/o history	-	-	-	19.55 \downarrow 2.34%
Ours w/o caption	-	-	-	20.06 \downarrow 1.11%
Ours w/o DTW	-	-	-	21.75 \downarrow 0.14%

4.9 Online Experiments

We test the proposed framework in real-world live streaming scenarios through online A/B testing. The experiment is conducted over four consecutive days, with traffic randomly assigned to either the baseline method or our method. In the baseline group, candidate live rooms are sorted by scores produced by a traditional recommendation model, e.g., Click-Through Rate s_{ctr} , Long-View-Through Rate s_{ltr} . Note that these scores focus on capturing the long-term relationship between streamers and users. In contrast, our method utilizes content-based CTR as an additional factor like s_{exp} . This term is capable of catching the highlight moment and show the most attractive live contents to users. The results show that our method achieve a 2.9% and 5.9% improvements in terms of CTR and live play duration, respectively, which demonstrate the

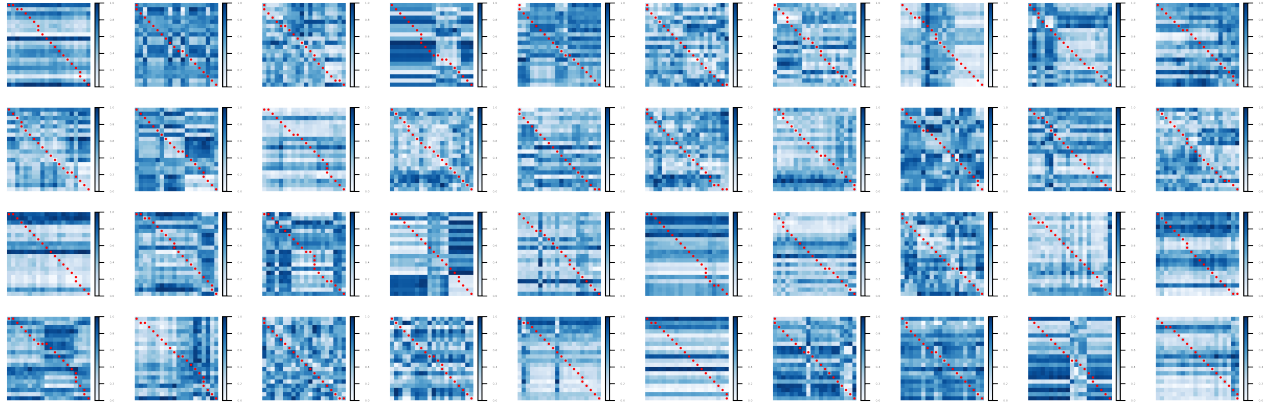


Figure 6: Visualization of Dynamic Time Warping (DTW) alignment results for the KLive dataset. Each subfigure represents the cosine similarity between video and text feature sequences. The red point denotes the DTW alignment path found by our algorithm.

effectiveness of our content based model.

$$\begin{aligned} s_{base} &= s_{ctr} + s_{lotr} + \dots \\ s_{exp} &= s_{ctr} + s_{lotr} + \dots + s_{ContentCTR}, \end{aligned} \quad (16)$$

5 CONCLUSION

In this paper, we study the task of Click-Through Rates prediction in live streaming scenario. We introduce ContentCTR, which utilizes a multimodal transformer to achieve frame-level CTR prediction. Specifically, we propose a streamer-personalized Perceiver Block that fuses ID embedding, visual, audio, and comment embedding. The decoder network outputs the final CTR prediction for each frame. To address the possible misalignment between video frames and texts, we carefully design a Dynamic Time Warping (DTW) alignment loss for optimization. Additionally, the boundary-aware constrained pairwise loss demonstrates better performance when combined with the pointwise loss. We conduct comprehensive experiments on both the KLive dataset and the public PHD dataset, which demonstrate the effectiveness of our methods in both streaming CTR prediction and video highlight detection tasks, compared with state-of-the-art methods. Moreover, the proposed method has been deployed online on the company’s short video platform and serves over three million daily users.

REFERENCES

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems* 35 (2022), 23716–23736.
- [2] Uttaran Bhattacharya, Gang Wu, Stefano Petrangeli, Viswanathan Swaminathan, and Dinesh Manocha. 2022. Show Me What I Like: Detecting User-Specific Video Highlights Using Content-Based Multi-Head Attention. In *Proceedings of the 30th ACM International Conference on Multimedia*. 591–600.
- [3] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to rank using gradient descent. In *Proceedings of the 22nd international conference on Machine learning*. 89–96.
- [4] Qiwei Chen, Huan Zhao, Wei Li, Pipei Huang, and Wenwu Ou. 2019. Behavior sequence transformer for e-commerce recommendation in alibaba. In *Proceedings of the 1st International Workshop on Deep Learning Practice for High-Dimensional Sparse Data*. 1–4.
- [5] Runnan Chen, Penghao Zhou, Wenzhe Wang, Nenglu Chen, Pai Peng, Xing Sun, and Wenping Wang. 2021. Pr-net: Preference reasoning for personalized video highlight detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7980–7989.
- [6] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [8] Avinash K Dixit. 1990. *Optimization in economic theory*. Oxford University Press, USA.
- [9] Ge Gao, Hongyan Liu, and Kang Zhao. 2023. Live streaming recommendations based on dynamic representation learning. *Decision Support Systems* 169 (2023), 113957.
- [10] Ana Garcia del Molino and Michael Gygli. 2018. Phd-gifs: personalized highlight detection for automatic gif creation. In *Proceedings of the 26th ACM international conference on Multimedia*. 600–608.
- [11] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. *arXiv preprint arXiv:1703.04247* (2017).
- [12] Li He, Hongxu Chen, Dingxian Wang, Shoaib Jameel, Philip Yu, and Guandong Xu. 2021. Click-through rate prediction with multi-modal hypergraphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 690–699.
- [13] Glenn Jocher. 2020. YOLOv5 by Ultralytics. <https://doi.org/10.5281/zenodo.3908559>
- [14] Maurice G Kendall. 1945. The treatment of ties in ranking problems. *Biometrika* 33, 3 (1945), 239–251.
- [15] Dohwan Ko, Joonmyung Choi, Juyeon Ko, Shinyeong Noh, Kyoung-Woon On, Eun-Sol Kim, and Hyunwoo J Kim. 2022. Video-text representation learning via differentiable weak temporal alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5016–5025.
- [16] Bin Liu, Chenxu Zhu, Guilin Li, Weinan Zhang, Jincai Lai, Ruiming Tang, Xiuqiang He, Zhenguo Li, and Yong Yu. 2020. Autofis: Automatic feature interaction selection in factorization models for click-through rate prediction. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2636–2645.
- [17] Tie-Yan Liu et al. 2009. Learning to rank for information retrieval. *Foundations and Trends® in Information Retrieval* 3, 3 (2009), 225–331.
- [18] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [19] Wantong Lu, Yantao Yu, Yongzhe Chang, Zhen Wang, Chenhui Li, and Bo Yuan. 2021. A dual input-aware factorization machine for CTR prediction. In *Proceedings of the twenty-ninth international conference on international joint conferences on artificial intelligence*. 3139–3145.
- [20] Meinard Müller. 2007. Dynamic time warping. *Information retrieval for music and motion* (2007), 69–84.
- [21] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).

- [22] Wentao Ouyang, Xiuwu Zhang, Li Li, Heng Zou, Xin Xing, Zhaojie Liu, and Yanlong Du. 2019. Deep spatio-temporal neural networks for click-through rate prediction. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2078–2086.
- [23] Qi Pi, Guorui Zhou, Yujing Zhang, Zhe Wang, Lejian Ren, Ying Fan, Xiaoqiang Zhu, and Kun Gai. 2020. Search-based user interest modeling with lifelong sequential behavior data for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2685–2692.
- [24] Jiarui Qin, Weinan Zhang, Xin Wu, Jiarui Jin, Yuchen Fang, and Yong Yu. 2020. User behavior retrieval for click-through rate prediction. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2347–2356.
- [25] Jérémie Rappaz, Julian McAuley, and Karl Aberer. 2021. Recommendation on live-streaming platforms: Dynamic availability and repeat consumption. In *Proceedings of the 15th ACM Conference on Recommender Systems*. 390–399.
- [26] Steffen Rendle. 2010. Factorization machines. In *2010 IEEE International conference on data mining*. IEEE, 995–1000.
- [27] Mrigank Rochan, Mahesh Kumar Krishna Reddy, Linwei Ye, and Yang Wang. 2020. Adaptive video highlight detection by learning from user history. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*. Springer, 261–278.
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [29] Hang Wang, Penghao Zhou, Chong Zhou, Zhao Zhang, and Xing Sun. 2022. PAC-Net: Highlight Your Video via History Preference Modeling. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIV*. Springer, 614–631.
- [30] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD'17*. 1–7.
- [31] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *Proceedings of the web conference 2021*. 1785–1797.
- [32] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metzger, Luke Zettlemoyer, and Christoph Feichtenhofer. 2021. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084* (2021).
- [33] Tzu-Wei Yang, Wen-Yuah Shih, Jiun-Long Huang, Wei-Chih Ting, and Pin-Chuan Liu. 2013. A hybrid preference-aware recommendation algorithm for live streaming channels. In *2013 Conference on Technologies and Applications of Artificial Intelligence*. IEEE, 188–193.
- [34] Sanshi Yu, Zhuoxuan Jiang, Dong-Dong Chen, Shanshan Feng, Dongsheng Li, Qi Liu, and Jinfeng Yi. 2021. Leveraging tripartite interaction information from live stream E-commerce for improving product recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3886–3894.
- [35] Shuai Zhang, Hongyan Liu, Jun He, Sanpu Han, and Xiaoyong Du. 2021. A deep bi-directional prediction model for live streaming recommendation. *Information Processing & Management* 58, 2 (2021), 102453.
- [36] Shuai Zhang, Hongyan Liu, Jun He, Sanpu Han, and Xiaoyong Du. 2021. Deep sequential model for anchor recommendation on live streaming platforms. *Big Data Mining and Analytics* 4, 3 (2021), 173–182.
- [37] Shuai Zhang, Hongyan Liu, Jun He, Sanpu Han, and Xiaoyong Du. 2021. Deep sequential model for anchor recommendation on live streaming platforms. *Big Data Mining and Analytics* 4, 3 (2021), 173–182.
- [38] Zhenhui Zhu, Zhi Yang, and Yafei Dai. 2017. Understanding the gift-sending interaction on live-streaming video websites. In *Social Computing and Social Media. Human Behavior: 9th International Conference, SCSM 2017, Held as Part of HCI International 2017, Vancouver, BC, Canada, July 9–14, 2017, Proceedings, Part I 9*. Springer, 274–285.