

Moment&Cross: Next-Generation Real-Time Cross-Domain CTR Prediction for Live-Streaming Recommendation at Kuaishou

Jiangxia Cao, Shen Wang, Yue Li, Shenghui Wang, Jian Tang, Shiyao Wang,
Shuang Yang, Zhaojie Liu, Guorui Zhou
Kuaishou Technology, Beijing, China
{caojiangxia, wangshen, liyue16, wangshenghui03, tangjian03, wangshiyao08,
yangshuang08, zhaotianxing, zhouguorui}@kuaishou.com

ABSTRACT

Kuaishou, is one of the largest short-video and live-streaming platform, compared with short-video recommendations, live-streaming recommendation is more complex because of: (1) temporarily-alive to distribution, (2) user may watch for a long time with feedback delay, (3) content is unpredictable and changes over time. Actually, even if a user is interested in the live-streaming author, it still may be an negative watching (e.g., short-view < 3s) since the real-time content is not attractive enough. Therefore, for live-streaming recommendation, there exists a challenging task: *how do we recommend the live-streaming at right moment for users?* Additionally, our platform’s major exposure content is short short-video, and the amount of exposed short-video is 9x more than exposed live-streaming. Thus users will leave more behaviors on short-videos, which leads to a serious data imbalance problem making the live-streaming data could not fully reflect user interests. In such case, there raises another challenging task: *how do we utilize users’ short-video behaviors to make live-streaming recommendation better?*

For the first challenge, we analyzed our data and observed an interesting phenomenon: when a live-streaming is at a high-light moment (e.g., Dancing), its click rate (CTR) will be increasing significantly. Inspired by our observation, we believe there exists a potential solution to automatically identify which live-streaming may be at its ‘high-light moment’, based on the wisdom of crowds of many user behaviors towards the live-streaming current moment. Therefore, our goal is to enable model to perceive all behaviors occurring as soon as possible, so that model can know which live-streaming is in the CTR increasing status. To achieve the idea, we have upgraded our data-streaming engine to real-time 30s report manner and devised a novel first-only mask learning strategy to supervise our model, named **Moment**. For the second challenge, we mainly follow the search-based interest modeling idea: first devise General Search Units (GSUs) to search users’ short-video/live-streaming history, and then use Extract Search Units (ESUs) to compress them. Besides, we also introduce a contrastive objective to align short-videos and live-streaming embedding spaces to enhance their correlation,

named **Cross**. We conduct extensive offline/online and ablation studies to verify our Moment&Cross effectiveness.

CCS CONCEPTS

• **Information systems** → **Recommender systems**.

KEYWORDS

Data-Streaming; Live-Streaming; Cross-Domain Recommendation

ACM Reference Format:

Jiangxia Cao, Shen Wang, Yue Li, Shenghui Wang, Jian Tang, Shiyao Wang., Shuang Yang, Zhaojie Liu, Guorui Zhou. 2024. Moment&Cross: Next-Generation Real-Time Cross-Domain CTR Prediction for Live-Streaming Recommendation at Kuaishou. In *Proceedings of ACM Conference (Conference’17)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nmmnnnnn.nmmnnnnn>

1 INTRODUCTION

Short-video and live-streaming platform like Kuaishou and Douyin have grown rapidly in recent years, attracting a lot of attention and accumulating a large number of active users. At Kuaishou, users always watch content at the **slide** page: the content on this page will automatically be played according to the user’s up and down scrolling on screen. Therefore, a powerful RecSys [9, 10] is the foundation of our service, to influence user watching experience to decide what content to watch next for users. Compared with the wide-explored short-video recommendation [29, 31, 34, 39], live-streaming recommendation [25] is a harder task since different natures of media as: (1) *Temporarily life-cycle*: Different from short-video has an eternal life-cycle to distribute, live-streaming is temporary (average 1 hour) to distribute in our system. (2) *Long-term feedback delay* [2]: Unlike short-videos, their average duration is about 55s, thus the report manner could quickly report all user behaviors to supervise model training. Nevertheless, the live-streaming is much longer, and some sparse and **valuable feedback may happen after watching half an hour**, e.g., users buy digital gifts for the live-streaming author. (3) *Dynamic content change* [30]: Different from short videos always playing from 0s, live-streaming content is ever-changing, thus for a live-streaming, a user watching at **different time points maybe have distinct behaviors**.

Therefore, our live-streaming RecSys needs to solve a very challenging problem: *how do we recommend the live-streamings at right moment for users?* To answer the question, we show two live-streaming author cases in Figure 1: (1) For the talent-show live-streaming: their authors spend a lot of time chatting with audiences, participating PK with other authors, and showing their talents sometimes. In terms of user behaviors, there is a significant difference between high-light moments and chatting: when the 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

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nmmnnnnn.nmmnnnnn>

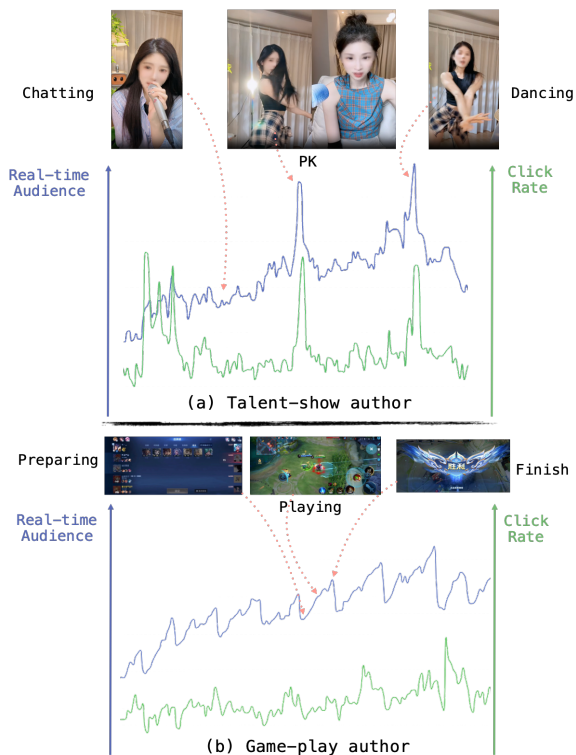


Figure 1: Typical live-streaming pattern between the ‘high-light moment’, real-time audience and the CTR trends.

shows dancing talent, the number of users entering live-streaming will increase significantly, and the user will exit after finish talent showing. (2) For the game-play live-streaming: their authors play game competitions one by one. During the author playing at one competition, the live-streaming real-time audience will continue to accumulate until the end of the competition. At the competition finished, the number of watching audience dropped significantly.

Actually, no matter what category the author is, users always tend to click the live-streaming for better watching live-streaming high-light moments, but the live-streaming content is ever-changing, so it is not an easy task to identify which live-streaming will show highlight moment. Fortunately, there may exist a potential solution to automatically identify which live-streaming is at its ‘high-light moment’, based on the wisdom of many user behaviors towards the live-streaming current moment: in Figure 1, the CTR trend shows highly consistency relation with user click behaviors (i.e., peaked and troughed at the same time). Therefore, if our model could **capture the increasing CTR trend**, then the model can discover potential highlight moments based on a large amount of user positive feedback.

Aside from the highlight moments capturing challenge, our live-streaming model has a more serious problem: data sparsity. On the slide page, users can watch short-videos and live-streamings in an interspersed manner according to the user’s up-and-down scrolling on screen. Nevertheless, the slide page is about 90% of the exposed content is short-videos, thus our live-streaming RecSys has the risk

could not fully learning the user interests to make precise CTR predictions. In such case, there raises another challenging task: *how do we utilize users’ rich short-video behaviors to make live-streaming recommendation better?* To answer the question, we first explain the workflow of our architecture (as shown in Figure 2). In industry, the different businesses are deployed separately, for instance, the users’ short-video real-time interactions (e.g., long-view, like and so on) are only assembled by the short-video data-streaming engine to be organized into specific training sample data formats. Then, the short-video model can consume short-video data-streaming to fit real-time data distribution to make precise recommendation. Since the different data-streaming engines will generate different data formats training samples, thus our live-streaming model can only be supervised by the users’ live-streaming data-streaming. Although we cannot consume short-video data-streaming directly, fortunately, we have constructed historical storage services to save users’ interaction logs [8, 13], and the data-streaming engine could send requests to obtain users’ interaction history from other businesses to **assemble them as a part of input features**. Therefore, we can align the live-streaming and short-video embedding space, then our model could utilize users short-video interests to determine which live-streaming with a similar style should be recommended.

In this paper, we present our effective and efficient solutions **Moment&Cross** - towards building the next-generation live streaming framework. For the first challenge, our goal is to encourage live-streaming model to be able to perceive what kind of live-streaming has the **CTR increasing trends**. Therefore, we need to utilize the real-time occurring user behaviors to train our model as soon as possible, to capture the real-time CTR status for each live-streaming. As shown in Figure 2, the CTR signal is first reported to our live-streaming data-streaming engine and then fed to our model. However, as many industrial RecSys implemented, the report module needs to wait for a while (e.g., 5 minutes) to collect enough behaviors and then report them to the data-streaming engine at once. Particularly, our live-streaming services also employ a fixed-window 5-minute style data-streaming to support our model for several years, however, a 5-minute feedback delay is not real-time enough to support our model to capture the CTR increasing trend. To this end, we have upgraded our training framework from Fast-Slow report manner to real-time 30s report manner, and devised a novel first-only mask learning strategy to supervise our model, named **Moment**. For the second challenge, our goal is to **exploit users historical short-video sequences** and align their embedding space with live-streaming. In fact, the user’s short-video history is too long to model directly [27] (e.g., a high-active user easily watch 10,000 short-videos in 1 month), thus we mainly follow the cascading search-based [24] interest modeling framework: (1) introduce General Search Units (GSUs) to retrospect user life-long history and then filter to obtain a top related item sequence (hundreds-level); (2) devise Exact Search Units (ESUs) to compress sequence information to obtain user interests, e.g, sequence pooling [15], target-item-attention [37]. Besides, we also introduce a contrastive [3] objective to align short-videos and live-streaming embedding spaces to enhance their correlation, thus our model could distribute live-streaming with a similar style according to users’ interests from rich short-video interaction history, named **Cross**. In summary, our contributions are as follows:

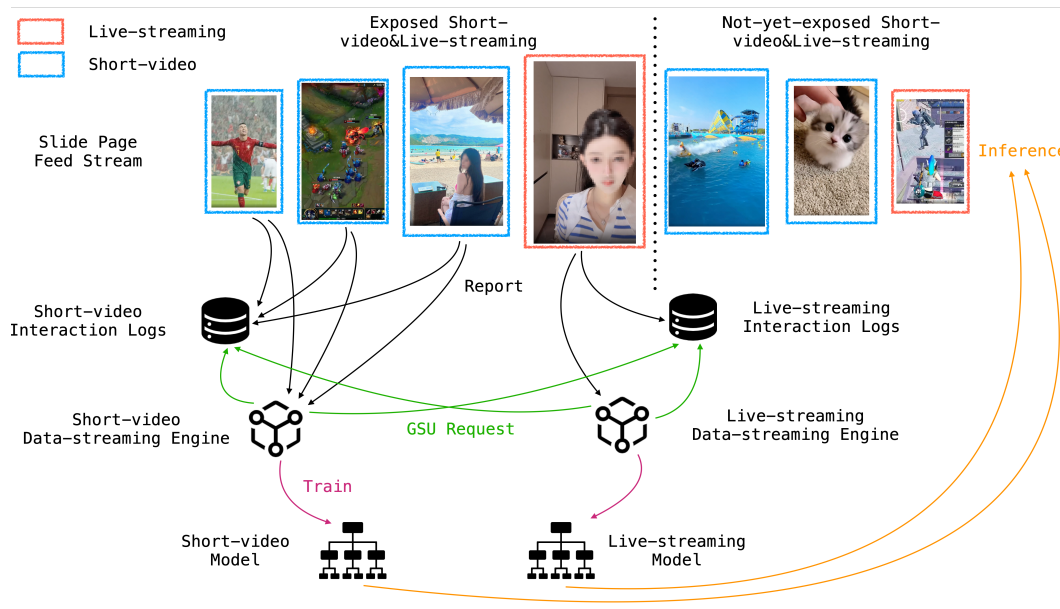


Figure 2: The Slide page RecSys architecture of Kuaishou short-video and live-streaming services, different services are separately with their own data-streaming and model. The only way to access users' other business logs, is by utilizing the 'interaction logs' storage services to retrospect historical user's short-video behaviors to find related small group items.

- We present a novel real-time learning framework to discover the 'high-light' live-streaming automatically, towards building next-generation live-streaming recommendation.
- We devise simple-yet-effective techniques to transfer user short-video interests for live-streaming recommendation.
- We conduct extensive offline and online experiments to validate our Moment&Cross, which has now been deployed on the live-streaming service at Kuaishou, serving 400 Million users.

2 MOMENT&CROSS AT KUAISHOU

The industry CTR prediction model [14] training process consists of two essential components: (1) a data-streaming engine to organize training samples features and labels, (2) a multi-task learning based model [23] to fit real training samples' interactions (e.g., click, like, long-view and others). In this Section, we first review the general preliminaries of our previous 5-min Fast-Slow live-streaming data-streaming engine building and CTR prediction model learning paradigm. We then show our novel real-time 30s data-streaming engine and its first-only mask learning strategy. Finally, we give our cross-domain techniques to capture the long-term and short-term short-video interaction patterns of users.

2.1 Preliminary: Fast-Slow 5-Min&1-Hour Data-Streaming

Data-streaming engines as the fundamental component of industry RecSys, a naive solution is to collect the user action logs to report after the item has been fully consumed, e.g., watch and swipe to next short-video, finish a listening song. Actually, the above solution is

'real-time' enough for short-video service, since user exit a short-video average within 1 minute, thus all interaction feedback will be collected in a short time. Nevertheless, in live-streaming service, users may watch for a very long time (e.g., 30min or 80min), if we still collect all actions when the user exits live-streaming, this will result in the model training not being 'real-time' enough.

Consequently, to our knowledge, many live-streaming data-streaming engines follow the **fixed-window** style to report and assemble training samples features and labels, e.g., 5-min report window. As one of the characteristics of live-streaming service, different behaviors occur with large different time distributions, and some valuable interactions are hard to observe within a small fixed-window, such as users gifting the author after watching for half an hour. To this end, we further extend the fixed-window style to the fast-slow windows for our live-streaming service to achieve a balance in Figure 3, i.e., **a fast-window reports all interactions for fast training, and a slow-window reports positive samples that have not been observed in the fast-window.**

Through statistics, we found that most users watch for less than 1 hour, thus we divided the user's watching experience into three periods to monitor user behaviors to supervise our model training:

- *5-minutes*: a small window to collect **observed all positive and negative behaviors as labels**, to obtain a fast data-streaming.
- *1-hour*: a large window only report the **missing positive labels** to correct the 5-minutes error as a slow data-streaming.
- *Ignored window*: this window **no longer report** any labels.

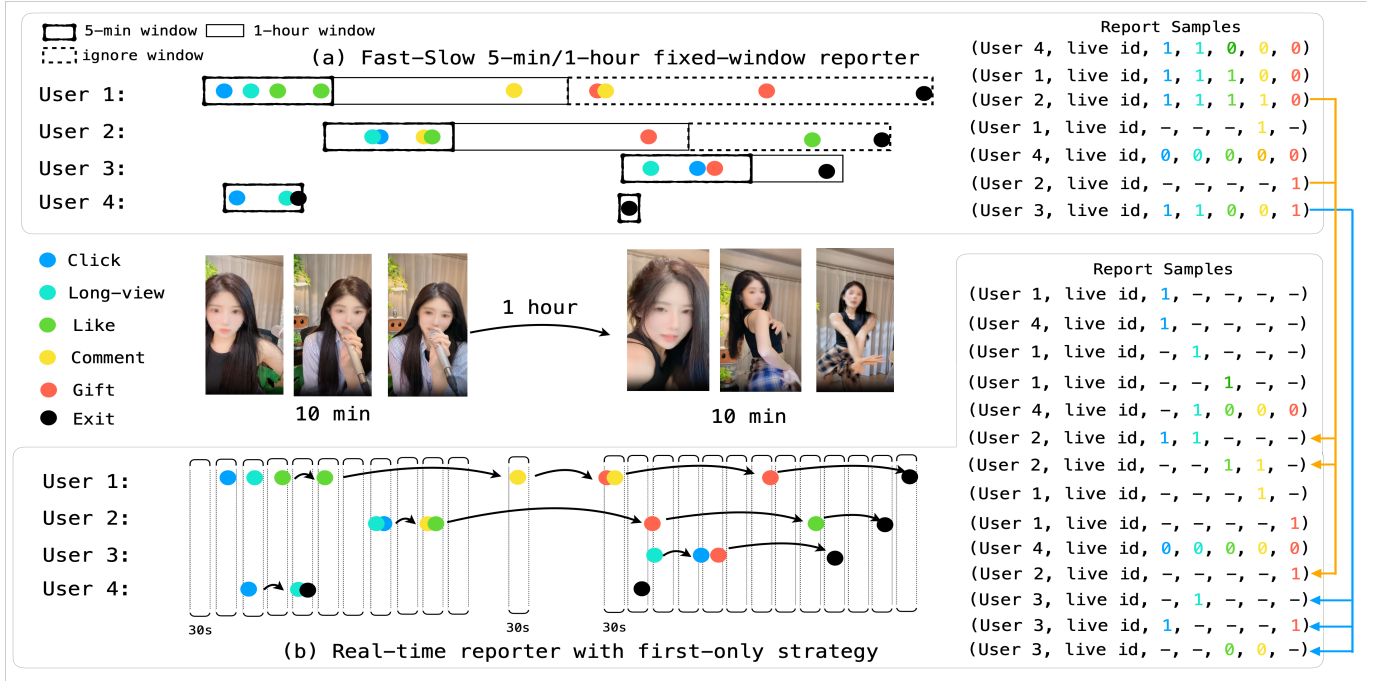


Figure 3: The report samples difference of produced training samples between fast-slow 5-min&1-hour data-streaming and real-time 30s data-streaming. We only show the simplest sample format (user, live-streaming, click, long-view, like, comment, gift). Specifically, for the fast-slow data-streaming, the fast flow reports 5-min window observed all user behaviors, the slow flow reports 5-min missing but 1-hour observed positive user behaviors. In our real-time data-streaming, we report users’ first positive behaviors immediately every 30 seconds and report all negative behavior when user exit a live-streaming. According to the report samples’ indicative relationship, our real-time data-streaming could produce training samples as soon as possible, to encourage model capturing the CTR increasing trends live-streaming.

2.2 Preliminary: CTR Model Training with Positive-Unlabeled Learning

In a broad sense, the CTR prediction [36] model is at the last point of RecSys [26], to rank the most related dozens of items for each user, also called fullrank model. Indeed, the fullrank [38] model is not only to predict the probability that a user would *click* the item candidate (i.e. CTR), and also predict the *long-view* probability (i.e. LVTR), *like* probability (i.e. LTR), *comment* probability (i.e. CMTR) and others XTRs at same time. According to these predicted probabilities, we can write some complex weighted calculations on these probabilities to control the final score to rank items.

Generally, the fullrank model learning process is formulated as a multi-task [33] binary classifications paradigm, where the goal is to learn a prediction function $f_{\theta}(\cdot)$ given the data-streaming training sample. For each sample consists of the user-item ID, features raw features V and several real labels to denote the action happens or not (i.e., $y^{ctr} \in \{0, 1\}$, $y^{ltr} \in \{0, 1\}$, $y^{cmtr} \in \{0, 1\}$ and others). Specifically, the raw features V mainly divide four types: user/item IDs, statistics/categories, historical interaction sequences and pre-trained LLM multi-modal [1, 4] embeddings, and these features are projected into a low-dimensional embedding as $V = [v_1, v_2, \dots, v_n]$, n indicates the number of features. In our live-streaming model, we hand-craft about $n > 400$ raw features

to represent user, item and context status. According to the input sample features and label, the model learning process is formed as:

$$\hat{y}^{ctr}, \hat{y}^{ltr}, \hat{y}^{cmtr}, \dots = f_{\theta}([v_1, v_2, \dots, v_n]) \quad (1)$$

where the $\hat{y}^{ctr}, \hat{y}^{ltr}, \hat{y}^{cmtr}, \dots$ are predicted probabilities, while the $f_{\theta}(\cdot)$ is a multi-task module can be implemented by MMoE [22], PLE [28]. Next, we utilize the users’ real behavior to supervise our model to optimize parameters. For the 5-min fast-flow samples, since they report all observed positive and negative labels, thus we train our model by minimizing the standard negative log-likelihood:

$$\mathcal{L}_{fast} = - \sum_{ctr, \dots} (y^{xtr} \log(\hat{y}^{xtr}) + (1 - y^{xtr}) \log(1 - \hat{y}^{xtr})) \quad (2)$$

For the 1-hour slow-flow samples, which only report the missing positive labels while masking other consistent positive labels, we employ the positive-unlabeled loss [11, 18, 19] to **retract the past "fake negative" error gradient** as:

$$\mathcal{L}_{slow} = - \sum_{missing} (\log(\hat{y}^{xtr}) - \log(1 - \hat{y}^{xtr})) \quad (3)$$

where *missing* represents the positive labels observed only in the 1-hour window. Through the collaboration of these two losses \mathcal{L}_{fast} and \mathcal{L}_{slow} , our model achieves a balance training between

Table 1: The 5-Min&1-Hour data streaming label consistency.

	Click	Long-View	Like	Comment	Gift	Gift Price
5-Min/1-Hour	96%	100%	79%	78%	68%	34%

effective and efficiency for live-streaming service. After the model convergence, we can push it as an online fullrank model to respond to real user requests and select the highest score items as follows:

$$\text{Ranking_Score} = (1 + \hat{y}^{ctr})^\alpha * (1 + \hat{y}^{ltr})^\beta * (1 + \hat{y}^{lcr})^\gamma * \dots \quad (4)$$

where α, β, γ are hyper-parameter to assemble all predicted probabilities as one ranking score to sort dozens of item candidates, here we only show a naive case of the weighted Ranking_Score.

2.3 Moment: Real-time 30s Data-Streaming with First-Only Label-Mask Learning

As our former version, the Fast-Slow 5-min&1-hour data-streaming and positive-unlabeled learning framework have been iterating for several years, which is a stable, reliable and proven learning framework designed for live-streaming. Although effective, it still has some drawbacks, for example: a smaller fixed window will inevitably miss some valuable positive labels, e.g., gift. We give the label consistency comparison between fast 5-min and sLow 1-hour data-streaming in Table 1, which describes the proportion of positive samples for major behavior (e.g., click, like, etc.). We can find that sparser behaviors have lower coverage of label consistency, especially the gift and gift price.

Furthermore, the fast 5-min window is still not real-time enough to support our model to capture the CTR increasing trend to solve the serious challenge: *how do we recommend the live-streaming at right moment for users?* To this end, we have upgraded our training framework from Fast-Slow report manner to real-time 30s report manner, to enable our model to perceive all behaviors occurring as soon as possible. Ideally, if a live-streaming at the ‘high-light moment’, **there will be a large number of positive gradients to optimize our model parameters in a short period, thus our model can know which live-streaming is in the CTR increasing status and then improve such live-streaming online CTR prediction scores to make it recommended for more users to watch.** However, the extremely small 30s window may bring some mismatched risks with our former data-streaming as:

- Fake negative: compared with 5-min fixed-window, if we also report all positive&negative user behaviors after watching 30s, will be introducing to lot of ‘fake negative’ labels since some behaviors are delayed rather than no-occurrence.
- Frequent report: compared the fast-slow 5-min&1-hour data-streaming are only report only once positive label for each behavior. our 30s real-time flow may be reported multiple times (e.g., a user could comment several times).
- Interaction beyond: the 30s real-time data-streaming may divide positive labels into several training samples in chronological order. In this way, the early positive behavior updated gradients may help predict the subsequent positive behaviors, for example:

the early click positive sample and long-view positive sample updated the model parameters two times, which may overestimate the predicted probability of latter like or comment behaviors.

To overcome the fake negative problem, inspired by the “**mask label**” idea of our former slow 1-hour flow, we introduce a reporting mechanism is that: the positive labels are reported immediately, while the rest negative labels are reported when user exits live-streaming. As a result, we found that although the report window is much smaller (i.e., 5-min \rightarrow 30s), the amount of data samples did not increase significantly (e.g, about 2 \times than Fast-Slow data-streaming), since the additional report samples mainly rely on the sparse behaviors interaction numbers, e.g., like, comment and gift. For the frequent report problem, we further introduce a **first-only mask strategy that we only learn the first positive occurrence for each behavior**, to align the learning rule with our former data-streaming. Therefore, our **Moment** first-only label-mask learning can be formed as:

$$\mathcal{L}_{moment} = -\sum_{first,exit} (y^{xtr} \log(\hat{y}^{xtr}) + (1 - y^{xtr}) \log(1 - \hat{y}^{xtr})) \quad (5)$$

where the *first* indicates the first positive labels of each behavior where others are masked, and the *exit* denotes the rest of behaviors’ positive or negative labels when user exit a live-streaming. Under the label mask setting, we can utilize Eq.(5) to replace Eq.(2) and Eq.(3) to support our model training and without the **long-term feedback delay** issue. For the interaction beyond risk, particularly, we did not observe this phenomenon, we assume the reason is that the model parameters optimizing is to **fit all users’ data distribution and is hard to overfit to a specific user live-streaming pattern.** The difference report between fast-slow and 30s real-time data-streaming is shown in Figure 3.

2.4 Cross: Short-Video Interest Transfer

On the slide page of our model deployment, the exposure content is about 90% are short-videos with 10% live-streaming, due to the uneven distribution of our traffic, we must consider the challenge: *how do we utilize users short-video behaviors to make live-streaming recommendation better?* As shown in Figure 2, the different business models are only allowed to consume their own training data-streaming, thus our live-streaming model can only be supervised by the users’ live-streaming behaviors. However, fortunately, we have constructed historical storage services to save users’ interaction logs, and the data-streaming engine could send requests to obtain users’ interaction history from other businesses to assemble them as a part of input features. In fact, in the ‘Interaction log’, we could retrospect the latest 10,000 watching item ID and access some side-information, e.g., time gap, content multi-modal tag, labels, etc. In order to model such a long sequence, a wide-used solution is the two cascading search-then-extract [9, 24, 27] idea: (1) introduce General Search Units (GSUs) to search user history and then filter to obtain a top hundreds related item sequence; (2) devise Exact Search Units (ESUs) to aggregate sequence information to compress users’ interests, e.g, sequence pooling, target-item-attention.

In our implementation, we introduce several GSUs to search different interactions from multi-aspects to find the related short-video with target live-streaming candidate, including:

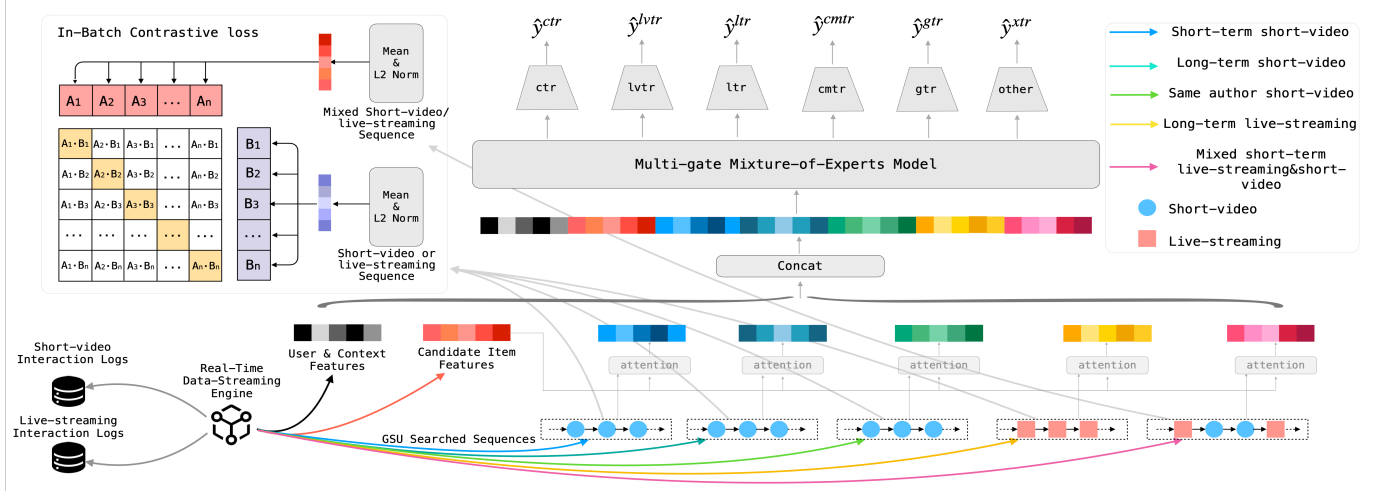


Figure 4: We introduce 5 short-video and live-streaming searched sequences to support our model: (1) we first conduct the contrastive mechanism to align them embedding space, (2) we then utilize the target attention mechanism to extract users' interests. (3) we finally concatenate the cross-domain short-video signal to predict each behavior probability.

- *Latest short-term short-video GSU* aims to search user's hundreds of latest short-video interactions V^{short} , which can reflect a precise user's short-term interest point.
- *Dot-product search long-term short-video GSU* search the short-video with the highest embedding similarity to the live-streaming candidate, denoted V^{long} , which could reflect whether users like this type of live-streaming.
- *Author ID hard search short-video GSU* aims to search the short-video history that with the same author ID, denoted $V^{aidhard}$, which can reflect precise user's interests of this author.
- *Dot-product search long-term live-streaming GSU* to obtain $V^{livelong}$, which could reflect whether users like this type of live-streaming according to similar short-video behaviors.
- *Long-view behavior latest mixed live-streaming&short-video GSU*, a positive long-view action hard search to obtain interested live-streaming and short-videos as mixed sequences V^{mixed} .

For notation brevity, we use the $V^{short} \in \mathbb{R}^{L \times D}$, $V^{long} \in \mathbb{R}^{L \times D}$, $V^{aidhard} \in \mathbb{R}^{L \times D}$, $V^{livelong} \in \mathbb{R}^{L \times D}$, $V^{mixed} \in \mathbb{R}^{L \times D}$ to represent different GSU sequences embeddings, where L is the sequence length. After obtaining sequence embeddings, we first conduct contrastive objectives to align their embedding spaces:

$$\begin{aligned}
 \mathcal{L}_{short}^{cl} &= \text{Contrastive}(\text{L2}(\text{Mean}(V^{mixed})), \text{L2}(\text{Mean}(V^{short}))) \\
 \mathcal{L}_{long}^{cl} &= \text{Contrastive}(\text{L2}(\text{Mean}(V^{mixed})), \text{L2}(\text{Mean}(V^{long}))) \\
 \mathcal{L}_{aidhard}^{cl} &= \text{Contrastive}(\text{L2}(\text{Mean}(V^{mixed})), \text{L2}(\text{Mean}(V^{aidhard}))) \\
 \mathcal{L}_{livelong}^{cl} &= \text{Contrastive}(\text{L2}(\text{Mean}(V^{mixed})), \text{L2}(\text{Mean}(V^{livelong})))
 \end{aligned} \quad (6)$$

where the $\text{Mean}(\cdot) : \mathbb{R}^{L \times D} \rightarrow \mathbb{R}^D$ is a simple pooling function to compress sequence representation, the $\text{L2}(\cdot)$ denote L2 normalization function, the $\text{Contrastive}(\cdot, \cdot)$ is a **in-batch sampling** function to gather negative samples to contrastive with. Inspired by C²DSR ??, we figure that the mixed live-streaming&short sequences can be the cornerstone to align with others since it has

some similarity with other sequences, but is not completely identical. Afterward, we conduct the target-item-attention as ESU module to achieve a fine-grained interests extraction according to target live-streaming candidate embedding V^{live} :

$$V_{ESU} = \text{target-item-attention}(V^{live}W^q, V^k, V^v), \quad (7)$$

where the V^{live} denotes all live-streaming side features of a training sample (e.g. item tags, Live ID, Author ID, etc.). After obtaining the enhanced cross-domain short-video interests, we concatenate them to estimate each interaction probability, as shown in Figure 4.

3 EXPERIMENTS

In this section, we conduct detailed offline experiments and online A/B test on our Kuaishou live-streaming services, to evaluate our proposed method Moment&Cross.

3.1 Base Models and Evaluation Metrics

As shown in Figure 4, in industry ranking model, the multi-gate mixture-of-experts model plays a vital role in estimating various interactions probabilities, and it has several choices, e.g., MMoE [22], CGC [28], PLE [28], AdaTT [20] and so on. In this paper, we select representative multi-task learning methods, CGC and PLE, to verify our solution effectiveness. We conduct extensive experiments to test our Moment&Cross ability at live-streaming services and we utilize two classic offline metrics to evaluate our ranking quality, AUC and GAUC [37] (user grouped AUC). After our model convergence, we further push it to our online A/B test platform to process real user requests at two applications, Kuaishou and Kuaishou Lite, and we report some major metrics to show our Moment&Cross improvements, e.g., Watch Time.

3.2 Overall Performance

Table 2 shows the performances of our Moment and Cross, respectively. Specifically, our online service needs to tackle billions of user

Table 2: Offline Moment&Cross results(%) in term AUC of GAUC on live-streaming services at Kuaishou.

Model Variants	Click		Effective-view		Long-view		Like		Comment		Gift	
	AUC	GAUC	AUC	GAUC	AUC	GAUC	AUC	GAUC	AUC	GAUC	AUC	GAUC
PLE (Moment&Cross)	82.75	65.57	78.99	64.96	85.08	75.71	91.48	74.67	92.78	74.74	96.32	74.71
CGC (Moment&Cross)	-0.18	-0.48	-0.13	-0.21	-0.23	-0.32	-0.13	-0.53	-0.06	-0.33	-0.01	-0.52
Cross w/o V^{short}	-0.90	-1.59	-0.25	-0.28	-0.51	-2.11	-1.23	-2.59	-0.92	-2.97	-0.86	-4.70
Cross w/o V^{long}	-0.14	-0.36	-0.19	-0.39	-0.04	-0.22	-0.13	-0.33	-0.19	-0.57	-0.05	-0.24
Cross w/o $V^{aidhard}$	-0.17	-0.40	-0.10	-0.19	-0.20	-0.31	-0.08	-0.17	-0.13	-0.26	-0.07	-0.50
Cross w/o $V^{livelong}$	-0.11	-0.25	-0.08	-0.11	-0.10	-0.11	-0.25	-0.55	-0.05	-0.16	-0.02	-0.15
Cross w/o V^{mixed} & \mathcal{L}^{cl}	-0.14	-0.40	-0.01	-0.20	-0.04	-0.54	-0.24	-0.42	+0.24	+0.12	+0.02	-0.79

Table 3: Online Moment&Cross A/B testing performance of live-streaming services at Kuaishou.

Applications	Modifications	Groups	Core Metrics			Interaction Metrics		
			Click	Watch Time	Gift Count	Like	Comment	Follow
Kuaishou	Moment	Total	+1.63%	+4.13%	-0.55%	-	-	+3.70%
		Total	+2.21%	+2.27%	+6.91%	+2.84%	+2.23%	+4.21%
	Corss	Low-Gift	+6.75%	+6.56%	+8.50%	+0.53%	+7.75%	+12.84%
		Middle-Gift	+4.27%	+3.38%	+9.21%	+3.77%	+6.47%	+7.08%
		High-Gift	+0.11%	+1.16%	+4.15%	+3.36%	+0.15%	+0.26%
Kuaishou Lite	Moment	Total	+0.64%	+1.85%	-1.22%	-	-	+2.79%
		Total	+2.72%	+2.48%	+8.91%	+1.05%	+4.58%	+4.37%
	Cross	Low-Gift	+2.94%	+5.69%	+24.75%	+5.75%	+3.94%	+7.07%
		Middle-Gift	+5.35%	+3.86%	+15.08%	+2.41%	+3.11%	+8.76%
		High-Gift	+0.52%	+1.93%	+4.42%	+0.21%	+4.18%	+1.37%

requests per day, and the improvement of 0.10% in offline evaluation of AUC and GAUC is significant enough to bring online gains. From the results, we have the following observations:

- (1) To investigate our real-time data-streaming effectiveness, we implement two multi-task variants: PLE(Moment&Cross) and CGC(Moment&Cross), where PLE is a double-layer stacked version of CGC, which is also the deployed model to support online request traffic. According Table 2, we can observe the PLE variant shows the expected confidence improvement compared to the CGC variant, which reveals our new real-time data-streaming could seamlessly support other models with the first-only label-mask learning strategy.
- To validate the Cross domain short-video interest effectiveness, we conduct ablation studies to test all GSU sequences effectiveness one by one, e.g., Cross without short-term short-video sequence. From Table 2, we have the following observations: (1) All of our Cross variants show significant performance degeneration, which indicates users' history short-video or live-streaming sequences could enhance our model to capture users interests more accurately. (2) Compared with the live-streaming based sequence, the short-video based sequences are more powerful in providing predictive information to empower live-streaming

ranking model, i.e., the short-term short-video sequence V^{short} could provide 0.9% improvement. The reason might be that users watching experience about 90% of the watching content are short-videos, leading their interest points will be better reflected in short-video historical sequence. which reveals that transferring the user cross-domain interest from rich short-video domain is a powerful technique to support our live-streaming service.

3.3 Online A/B Test

To quantify the contribution of Moment&Cross could brings to our live-streaming services, we push the corresponding modification to our online A/B test system to serve as a ranking model at two applications, the Kuaishou and Kuaishou Lite. We evaluate model performance based on the core metrics and interaction metrics, e.g., Watch-Time, Gift Price, Click, etc. Table 3 reports our online results of Moment and Cross individually, here we further shows the fine-grained improvement value of three user groups to test our short-video interests transferring: Low/Middle/High-Gift users and total users. Specifically, since a long time has passed of our Moment modification, we unfortunately lost the likes and comments interaction metric results, but our core metrics were retained. According Table 3, we can find real-time Moment data-streaming

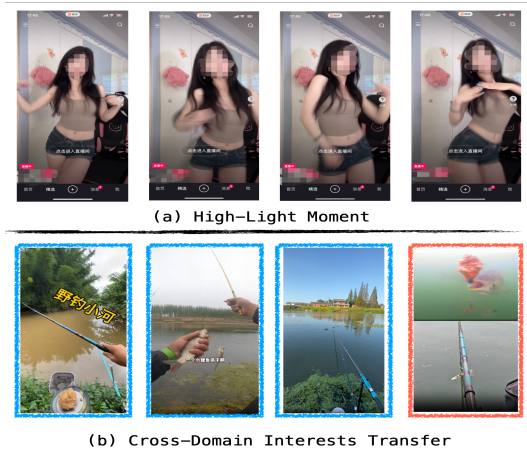


Figure 5: (a) Moment could enhance our model to perceive the high-light live-streamings; (b) Short-video cross-domain interests could help our system find related live-streamings.

trained model achieves a large improvement at Click +1.63%/+0.64% and Watch Time +4.13%/+1.85%, which indicates that accelerating the model real-time training efficiency is crucial for live-streaming recommendation. Regarding the slightly negative -0.55%/-1.22% of Gift Count, this is because gift is an **unstable metric** and within a reasonable range in our system. Our Cross achieves a large improvement of 2.27%/2.48% at watch-time and 6.91%/8.91% at gift price, the results demonstrate that our Cross could contribute to our system significantly. Further, the Low-Gift user group shows the biggest growth than others, which indicates the rich short-video signal is helpful for our system to alleviate data-sparse problem.

3.4 Case Study

This Section discovers the certain experience influence of Moment&Cross, we give three cases to demonstrate its effectiveness:

- As shown in Figure 5(a), we find out that the Slide page could recommend more ‘high-light moment’ of talent-show authors, which indicates the 30s real-time data-streaming with mask learning paradigm could capture the CTR increasing trends to detect ‘high-light moment’ accurately. And such phenomenon will not only bring a better experience to user but also allow the efforts of our platform’s authors to be seen by more users, building a better environment for live-streaming.
- For the cross-domain interests transferring, as shown in Figure 5(b), we observe that our system could feed more related live-streamings for users, e.g., a user long-viewed some fishing videos and then recommend outdoor fishing live-streamings. To be specific, fishing live-streaming is a minority category, without the help of short-video signals, it is hard for our system to successfully recommend this live-streaming.
- Moreover, since Kuaishou is a mainly short-video platform, many users do not have the habit of watching live-streaming. We are also curious whether our model can encourage more users to watch live-streaming with the help of short video signals. Here

Table 4: Short-video interest effects of different user group.

Metrics	Live Activate User Groups			
	Low	Middle	High	Full
Click	-1.45%	-1.35%	-1.24%	-0.95%
Long-view	-1.79%	-1.74%	-1.74%	-1.57%
Like	-1.88%	-1.97%	-1.86%	-1.25%
Comment	-1.48%	-1.67%	-1.64%	-1.14%
Gift	-1.47%	-1.67%	-1.53%	-0.94%

we divide our users into four types (e.g., low/middle/high/full-activate user groups) to show our model’s predictive ability. From Table 4, we can find the low-activate group shows the most significant improvements than other groups, which indicates our cross-domain interest transfer could effectively discover potential user group for live-streaming service.

4 RELATED WORKS

In recent years, live-streaming has become a fashionable phenomenon, with a large number of professional authors relying on live-streaming media to interact with the audience. Different from other recommendation scenarios that connect users with some items, the live-streaming aims to link users with their interested author, thus the pioneer work LiveRec [25] considering the user-author repeated consumption relationship with self-attention mechanism. Further, users could watch a live-streaming for a long time, the [12] devise a loss re-weighted strategy, which adjusts loss by differing amounts of minutes watched. To consider the live-streaming multi-modal information effects, the MTA [30] and ContentCTR [16] further introduce some multi-modal components to fuse textual, image frame information, and the Sliver [21] introduces a re-recommendation mechanism to capture dynamic live-streaming change. Besides, the recent aims to capture the user-author and author-author relation, MMBee [17] provides a graph representation learning and meta-path based behavior expansion strategy to enrich user and item multi-hop neighboring information. In addition, the recent progress also points out the cross-domain [5–7] signal ability for live-streaming, the DIAGAE [35] utilizes live-streaming domain user representations to align with other rich services learned user representations. The eLiveRec [32] aims to improve e-commerce live streaming recommendation, which devises a disentangled encoder to learn user’s live-streaming and product shared intentions and live-streaming specific intentions. Compared with them, our Moment&Cross has the distinct motivations that we aim to solve the following two problems: (1) *how do we recommend the live-streamings at right moment for users?* (2) *how do we utilize users rich short-video behaviors to make live-streaming recommendation better?* Further, the used solutions are totally different: (1) we upgrade our data-streaming and devise a novel first-only mask learning strategy, (2) we introduce search-based framework to exploit rich domain interaction sequences with a contrastive objective.

5 CONCLUSIONS

In this paper, we propose the Moment&Cross, towards to build the next-generation recommendation model for live-streaming services. Specifically, we first explain the background of live-streaming at Kuaishou, and then present two questions: (1) *how do we recommend the live-streamings at right moment for users?* (2) *how do we utilize users rich short-video behaviors to make live-streaming recommendation better?* For the first challenge, we blame the reason that 5-min fixed report window is hard to perceive live-streaming hot trends, and we describe a real-time 30s report mechanism with first-only mask learning paradigm to alleviate the problem. For the second challenge, we first devise several GSUs to search short-video and live-streaming sequences from user historical logs, and then introduce a contrastive objective to align them with representation space to support our ranking model. Extensive offline and online experimental results on our industrial real-time 30s data-streaming demonstrate the effectiveness of Moment&Cross at live-streaming services. Further, detailed analyses from various perspectives show the effectiveness of Moment and Cross, respectively.

REFERENCES

- [1] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *Arxiv* (2023).
- [2] Ashwinkumar Badanidiyuru, Andrew Evdokimov, Vinodh Krishnan, Pan Li, Wynn Vonnegut, and Jayden Wang. 2021. Handling many conversions per click in modeling delayed feedback. *Arxiv* (2021).
- [3] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *International Conference on Machine Learning (ICML)*.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Conference on Neural Information Processing Systems (NeurIPS)*.
- [5] Jiangxia Cao, Shaoshuai Li, Bowen Yu, Xiaobo Guo, Tingwen Liu, and Bin Wang. 2023. Towards universal cross-domain recommendation. In *ACM International Conference on Web Search and Data Mining (WSDM)*.
- [6] Jiangxia Cao, Xixun Lin, Xin Cong, Jing Ya, Tingwen Liu, and Bin Wang. 2022. Disencdr: Learning disentangled representations for cross-domain recommendation. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*.
- [7] Jiangxia Cao, Jiawei Sheng, Xin Cong, Tingwen Liu, and Bin Wang. 2022. Cross-domain recommendation to cold-start users via variational information bottleneck. In *IEEE International Conference on Data Engineering (ICDE)*.
- [8] Yue Cao, Xiaojiang Zhou, Jiaqi Feng, Peihao Huang, Yao Xiao, Dayao Chen, and Sheng Chen. 2022. Sampling is all you need on modeling long-term user behaviors for CTR prediction. In *ACM International Conference on Information and Knowledge Management (CIKM)*.
- [9] Jianxin Chang, Chenbin Zhang, Zhiyi Fu, Xiaoxue Zang, Lin Guan, Jing Lu, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, et al. 2023. TWIN: Two-stage interest network for lifelong user behavior modeling in CTR prediction at kuaishou. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [10] Jianxin Chang, Chenbin Zhang, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, and Kun Gai. 2023. Pepnet: Parameter and embedding personalized network for infusing with personalized prior information. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [11] Olivier Chapelle. 2014. Modeling delayed feedback in display advertising. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [12] Edgar Chen, Mark Ally, Eder Santana, and Saad Ali. 2022. Weighing dynamic availability and consumption for Twitch recommendations. *Arxiv* (2022).
- [13] Qiwei Chen, Changhua Pei, Shanshan Lv, Chao Li, Junfeng Ge, and Wenwu Ou. 2021. End-to-end user behavior retrieval in click-through rate prediction model. *Arxiv* (2021).
- [14] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhya, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *ACM Conference on Recommender Systems (RecSys Workshop)*.
- [15] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *ACM Conference on Recommender Systems (RecSys)*.
- [16] 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. *Arxiv* (2023).
- [17] Jiaxin Deng, Shiyao Wang, Yuchen Wang, Jiansong Qi, Liqin Zhao, Guorui Zhou, and Gaofeng Meng. 2024. MMBee: Live Streaming Gift-Sending Recommendations via Multi-Modal Fusion and Behaviour Expansion. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [18] Zhigang Huangfu, Gong-Duo Zhang, Zhengwei Wu, Qintong Wu, Zhiqiang Zhang, Lihong Gu, Jun Zhou, and Jinjie Gu. 2022. A multi-task learning approach for delayed feedback modeling. In *International World Wide Web Conference Companion*.
- [19] Sofia Ira Ktena, Alykhan Tejani, Lucas Theis, Pranay Kumar Myana, Deepak Dilipkumar, Ferenc Huszár, Steven Yoo, and Wenzhe Shi. 2019. Addressing delayed feedback for continuous training with neural networks in CTR prediction. In *ACM Conference on Recommender Systems (RecSys)*.
- [20] Danwei Li, Zhengyu Zhang, Siyang Yuan, Mingze Gao, Weilin Zhang, Chaofei Yang, Xi Liu, and Jiyan Yang. 2023. AdaTT: Adaptive Task-to-Task Fusion Network for Multitask Learning in Recommendations. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [21] Fengqi Liang, Baigong Zheng, Liqin Zhao, Guorui Zhou, Qian Wang, and Yanan Niu. 2024. Ensure Timeliness and Accuracy: A Novel Sliding Window Data Stream Paradigm for Live Streaming Recommendation. *Arxiv* (2024).
- [22] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. 2018. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [23] Xiao Ma, Liqin Zhao, Guan Huang, Zhi Wang, Zelin Hu, Xiaoqiang Zhu, and Kun Gai. 2018. Entire space multi-task model: An effective approach for estimating post-click conversion rate. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*.
- [24] 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 *ACM International Conference on Information and Knowledge Management (CIKM)*.
- [25] Jérémie Rappaz, Julian McAuley, and Karl Aberer. 2021. Recommendation on live-streaming platforms: Dynamic availability and repeat consumption. In *ACM Conference on Recommender Systems (RecSys)*.
- [26] Steffen Rendle. 2010. Factorization machines. In *IEEE International Conference on Data Mining (ICDM)*.
- [27] Zihua Si, Lin Guan, ZhongXiang Sun, Xiaoxue Zang, Jing Lu, Yiqun Hui, Xingchao Cao, Zeyu Yang, Yichen Zheng, Dewei Leng, et al. 2024. TWIN V2: Scaling Ultra-Long User Behavior Sequence Modeling for Enhanced CTR Prediction at Kuaishou. *Arxiv* (2024).
- [28] Hongyan Tang, Junjing Liu, Ming Zhao, and Xudong Gong. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning (MTL) Model for Personalized Recommendations. In *ACM Conference on Recommender Systems (RecSys)*.
- [29] Shisong Tang, Qing Li, Dingmin Wang, Ci Gao, Wentao Xiao, Dan Zhao, Yong Jiang, Qian Ma, and Aoyang Zhang. 2023. Counterfactual video recommendation for duration debiasing. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- [30] Dinghao Xi, Liumin Tang, Runyu Chen, and Wei Xu. 2023. A multimodal time-series method for gifting prediction in live streaming platforms. *Information Processing & Management (IPM)* (2023).
- [31] Jing Yan, Liu Jiang, Jianfei Cui, Zhichen Zhao, Xingyan Bin, Feng Zhang, and Zuotao Liu. 2024. Trinity: Syncretizing Multi-/Long-tail/Long-term Interests All in One. *Arxiv* (2024).
- [32] Yixin Zhang, Yong Liu, Hao Xiong, Yi Liu, Fuqiang Yu, Wei He, Yonghui Xu, Lizhen Cui, and Chunyan Miao. 2023. Cross-domain disentangled learning for e-commerce live streaming recommendation. In *IEEE International Conference on Data Engineering (ICDE)*.
- [33] Yu Zhang and Qiang Yang. 2022. A Survey on Multi-Task Learning. *IEEE Transactions on Knowledge and Data Engineering (TKDE)* (2022).
- [34] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumbhakar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending what video to watch next: a multitask ranking system. In *ACM Conference on Recommender Systems (RecSys)*.
- [35] Jiawei Zheng, Hao Gu, Chonggang Song, Dandan Lin, Lingling Yi, and Chuan Chen. 2023. Dual Interests-Aligned Graph Auto-Encoders for Cross-domain Recommendation in WeChat. In *ACM International Conference on Information and Knowledge Management (CIKM)*.
- [36] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *AAAI Conference on Artificial Intelligence (AAAI)*.
- [37] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep Interest Network for Click-Through Rate Prediction. In *ACM SIGKDD Conference on Knowledge Discovery*

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

Jiangxia Cao, Shen Wang, Yue Li, Shenghui Wang, Jian Tang, Shiyao Wang,
Shuang Yang, Zhaojie Liu, Guorui Zhou

and Data Mining (KDD).

[38] Jieming Zhu, Jinyang Liu, Shuai Yang, Qi Zhang, and Xiuqiang He. 2021. Open benchmarking for click-through rate prediction. In *ACM International Conference on Information and Knowledge Management (CIKM)*.

[39] Yongchun Zhu, Jingwu Chen, Ling Chen, Yitan Li, Feng Zhang, and Zuotao Liu. 2024. Interest Clock: Time Perception in Real-Time Streaming Recommendation System. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*.