

Hydrus: Improving Personalized Quality of Experience in Short-form Video Services

Zhiyu Yuan*

School of Software, Tsinghua University
Kuaishou
Beijing, China
yzy18@tsinghua.org.cn

Gang Wang

Kuaishou
Beijing, China
wanggang@kuaishou.com

Kai Ren*

Kuaishou
Beijing, China
kair@alumni.cmu.edu

Xin Miao[†]

School of Software, Tsinghua University
Beijing, China
miaoxin@tsinghua.edu.cn

ABSTRACT

Traditional approaches to improving users' quality of experience (QoE) focus on minimizing the latency on the server side. Through an analysis of 15 million users, however, we find that for short-form video apps, user experience depends on both response latency and recommendation accuracy. This observation brings a dilemma to service providers since improving recommendation accuracy requires adopting complex strategies that demand heavy computation, which substantially increases response latency.

Our motivation is that users' sensitivity to response latency and recommendation accuracy varies greatly. In other words, some users would accept a 20ms increase in latency to enjoy higher-quality videos, while others prioritize minimizing lag above all else. Inspired by this, we present Hydrus, a novel resource allocation system that delivers the best possible personalized QoE by making tradeoffs between response latency and recommendation accuracy. Specifically, we formulate the resource allocation problem as a utility maximization problem, and Hydrus is guaranteed to solve the problem within a few milliseconds.

We demonstrate the effectiveness of Hydrus through offline simulation and online experiments in Kuaishou, a massively popular video app with hundreds of millions of users worldwide. The results show that Hydrus can increase QoE by 35.6% with the same latency or reduce the latency by 10.1% with the same QoE. Furthermore, Hydrus can achieve 54.5% higher throughput without a decrease in QoE. In online A/B testing, Hydrus significantly improves click-through rate (CTR) and watch time; it can also reduce system resource costs without sacrificing QoE.

CCS CONCEPTS

• **Information systems** → *Web services*; **Personalization**.

*Both authors contributed equally to this research.

[†]Xin Miao is the corresponding author.



This work is licensed under a Creative Commons Attribution International 4.0 License.

SIGIR '23, July 23–27, 2023, Taipei, Taiwan

© 2023 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9408-6/23/07.

<https://doi.org/10.1145/3539618.3591696>

KEYWORDS

personalized quality of experience, resource allocation, short-form video services

ACM Reference Format:

Zhiyu Yuan, Kai Ren, Gang Wang, and Xin Miao. 2023. Hydrus: Improving Personalized Quality of Experience in Short-form Video Services. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23)*, July 23–27, 2023, Taipei, Taiwan. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3539618.3591696>

1 INTRODUCTION

Large-scale Internet service providers have long struggled with improving end-to-end performance while accommodating a growing user base. This is particularly challenging for short-form video services, which require high throughput and low latency to deliver a satisfying user experience. For example, the average daily active users (DAUs) on the Kuaishou app in 2022 were 355.7 million, representing an increase of 15.4% from 308.2 million in 2021 and 34.4% from 264.6 million in 2020 [19]. Clearly, the explosive growth of users brings an unprecedented challenge in improving user QoE.

Prior work [1, 5, 9, 24, 33] demonstrates that the relationship between latency and QoE is sigmoidal in traditional Internet services (e.g., web search), as shown in Figure 1(a). This is intuitive because users will have a better experience if a web page takes less time to load. Therefore, a rich literature has developed around cutting tail/average latency to enhance user QoE in various ways, such as addressing the task size [12], queue time [34], cache [4, 20, 21] and replica selection [25]. However, we found that reducing latency may not be an effective way to improve user QoE in short-form video services, based on our in-depth analysis of 15 million Kuaishou users. This is due to two main observations:

- **novel relationship:** we observed a quadratic-like curve between response latency¹ and QoE in short-form video services, which is not consistent with the typical sigmoid-like curve in traditional services. As shown in Figure 1(b), using more complex recommendation strategies will result in an increase in both response latency and recommendation accuracy (the gradually darker blue color). In Phase A, despite the increase in latency, QoE continues

¹In our work, response latency refers only to the time required for backend computing tasks, and we do not consider external factors such as transmission latency.

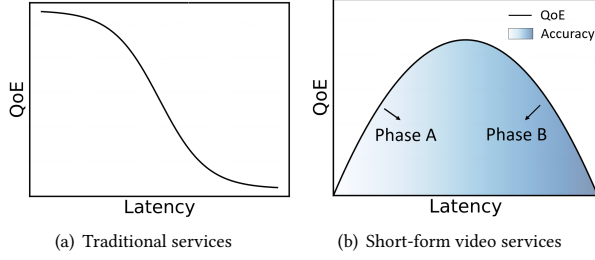


Figure 1: Comparison of the QoE-Latency relationship between traditional Internet services and short-form video services. (a) sigmoid-like curve, (b) quadratic-like curve.

to improve as the recommendation accuracy dominates. However, once the response latency exceeds the user’s tolerance, QoE gradually decreases in Phase B, even if higher-quality videos are accessible.

- **personalized QoE:** we observed that the QoE sensitivity to response latency and recommendation accuracy varies among these users, with users experiencing a peak in QoE at different latency levels. Furthermore, the magnitude of the QoE peak also differs among users. This variability in QoE sensitivity may be due to individual differences in preferences, expectations for recommendation accuracy, and tolerance for response latency. We elaborate on this observation in § 2.3.

At a high level, the novel relationship originates from the fundamentally conflicting trend in user expectations of lower latency and more accurate recommendations. Embracing the reality that we cannot satisfy the desire of all users with limited resources, we should leverage the different QoE sensitivity among users to personalized allocate resources instead of blindly reducing latency or improving recommendation accuracy. Therefore, we propose a novel effective solution to improve user QoE by personalized allocating resources according to different QoE sensitivities between response latency and recommendation accuracy.

The critical challenges associated with the solution are as follows. The first challenge is to model user heterogeneity. It’s a difficult and resource-intensive task to model hundreds of millions of users’ QoE sensitivity to response latency and recommendation accuracy. The second challenge is to make accurate real-time allocation decisions that improve overall QoE. It requires an extremely high-efficiency of resource allocation algorithm since an online short-form video service can receive tens of thousands of requests per second. The third challenge is that the algorithm should be sufficiently intelligent to make full use of the available resources because the number of free resources changes throughout the day.

To address these challenges, we present Hydrus, a novel resource allocation system to improve overall QoE by making personalized tradeoffs between response latency and recommendation accuracy according to user heterogeneity. Hydrus executes two steps to complete the resource allocation for a single request. First, Hydrus employs a parameter solver module to run a highly efficient algorithm that solves the strategy allocation problem and saves the personalized parameters in a database. The module runs offline,

ensuring it doesn’t create any additional latency on the online services. In the second step, Hydrus utilizes a strategy decision module to accurately and succinctly estimate the user QoE and resource cost of each strategy. Leveraging the previously solved parameters, the system then determines the best strategy allocation for each request to improve the overall QoE. Finally, it publishes the log information to the message queue for the parameter solver module to use in future calculations. We demonstrate the effectiveness and efficiency of Hydrus through trace-driven experiments and online evaluation in Kuaishou’s large-scale production system.

In summary, the paper makes the following main contributions:

- We found a novel quadratic-like relationship between QoE and response latency in short-form video services. Therefore, we propose to personalized allocate resources according to different QoE sensitivities between response latency and recommendation accuracy.
- We implement Hydrus, a novel resource allocation system, to improve QoE by making tradeoffs between response latency and recommendation accuracy according to user heterogeneity. Hydrus takes two steps to determine the best strategy allocation for each request.
- We evaluated the effectiveness and efficiency of Hydrus through offline trace-driven simulation and online experiments in Kuaishou. The offline simulation showed that Hydrus achieved a 35.6% higher QoE or 10.1% reduction in latency compared to the baseline. Furthermore, Hydrus is able to serve 54.5% more concurrent requests. In online experiments, Hydrus significantly improved CTR and watch time; it can also reduce system resource costs without sacrificing QoE.

The rest of the paper is organized as follows. We discuss our motivation in detail in § 2. The problem formulation and system design are described in § 3 and § 4, respectively. The evaluation of Hydrus and results analysis are presented in § 5. The related work is discussed in § 6. Finally, we conclude in § 7.

2 MOTIVATION

We collect millions of user data from the large-scale short-form video recommendation system of Kuaishou. Then, we present the analysis results to support our motivation. Finally, we show the potential QoE improvement achieved by making tradeoffs between response latency and recommendation accuracy in a simulation experiment.

2.1 Dataset

The dataset consists of user information and request traces served by the short-form video recommendation system from July 12 to July 15, 2022. It includes approximately 15M unique users and 428M traces, as summarized in Table 1. We record the following information for every request trace: user profile, recommendation strategy, response latency, and video completion rate (VCR).

- *user profile:* it is a collection of settings and information associated with a user that includes critical information used to identify an individual, such as user ID, gender, and age. These data are attached to each request and are used to analyze user heterogeneity in § 2.2.

	Hot Page	Nearby Page	Follow Page
User	9.5M	4.3M	1.2M
Trace	276.8M	123.7M	27.6M

Table 1: The size of the collected dataset in different pages from Kuaishou.

- *trace data* is the primary source data for our motivation. In detail, the dataset records the recommendation strategy, response latency of each request, and the corresponding user QoE.
 - *recommendation strategy*: short-form video services employ a variety of strategies, and the differences among strategies may be reflected in different recommendation models, numbers of service nodes, numbers of candidate videos, and so forth. For example, in the ranking stage, the choice of ranking 20 or 30 videos from the previous stage represents two distinct strategies. In this dataset, we collect the number of service nodes as different recommendation strategies. It is easy to understand that more complex strategies lead to an increase in the number of service nodes, which in turn improves the accuracy of recommendations.
 - *response latency*: a large-scale recommendation system often treats the sequence of different jobs and their dependencies as directed acyclic graphs (DAGs) [7, 22], such as fetching user profiles from the user database and ranking videos based on scores. The response latency represents the time to process all the jobs in the DAG, and we measure the latency using the timestamps tagged in each job. As described in § 1, our work does not consider the transmission latency (e.g., from server to user).
 - *video completion rate*: we choose VCR to estimate user QoE, which refers to the percentage of viewers watching a video from start to finish. We prefer to use VCR rather than watch duration [32] or user engagement [33] as the measure of QoE because VCR can eliminate the effects of different video durations. Furthermore, it can perfectly combine response latency and recommendation accuracy to characterize the service quality for the request. For example, VCR will decrease due to a long waiting time or unattractive videos.

2.2 Novel QoE-Latency relationship

Regarding traditional Internet services (e.g., web search), some work has measured user QoE by the difference between the start and end timestamps of the web session. In this way, they obtained a sigmoid-like relationship between QoE and latency. However, the measured result is different in short-form video services.

A request in short-form video services is handled as follows. When a request reaches the backend, the services fetch the user profile and return 10 to 15 personalized videos to the front end of the mobile application. Finally, the feedback of these videos (e.g., like, comment, watch time) is sent to the server by another request. Finally, we record the watch time to generate the corresponding VCR as a QoE value for every request.

We used the collected data over four days to explore the relationship between user QoE and latency. We first observed that an increase in the number of service nodes leads to longer response

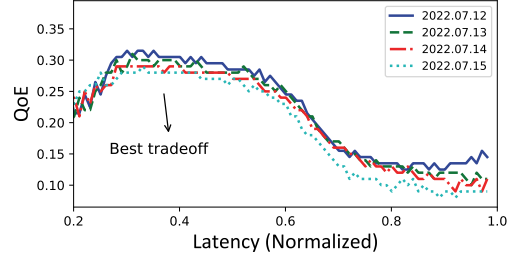


Figure 2: We collect and analyze trace data from July 12 to 15, 2022, observing a consistent and novel quadratic-like relationship between QoE and latency. Therefore, reducing latency alone will not lead to improved QoE.

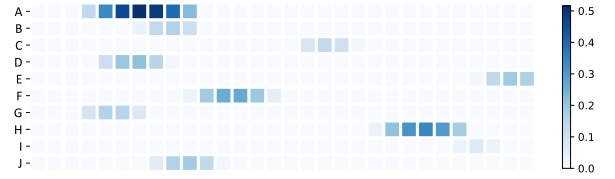


Figure 3: We visualize the QoE heatmap by randomly choosing ten users and calculating their average QoE at different latencies. Darker color represents a higher QoE value, which shows QoE sensitivity varies greatly among users.

times. Then we obtain four quadratic-like curves, which illustrate that QoE initially increases but then sharply drops as latency increases, as shown in Figure 2. Each point on the curve represents the average VCR at a given latency. This finding suggests that more complex recommendation strategies may not always result in improved QoE, as response latency tends to increase accordingly.

The novel relationship reveals that the developers of short-form video services often allocate significant resources to improve recommendation accuracy but overlook the impact of increasing latency on user experience. It is therefore wise to improve QoE by making tradeoffs between response latency and recommendation accuracy. Ideally, we should find the best tradeoff for each user, as shown in Figure 2, to maximize QoE with the least amount of resources.

2.3 User heterogeneity

The novel QoE-Latency relationship inspires us to make tradeoffs between response latency and recommendation accuracy to improve QoE, rather than simply reducing latency. However, this approach is challenging to implement because of user heterogeneity in QoE sensitivity. We obtained the historical trace data for the randomly selected users and visualized their average QoE using a heatmap, as shown in Figure 3. Although we only plotted the QoE data of 10 users, we observed a similar pattern in the QoE heatmap for other users in the dataset. It reveals that QoE sensitivity to response latency and recommendation accuracy varies among users. This implies that the same amount of system resources can generate different QoE benefits for different users. Consequently, we can broadly classify users into three categories:

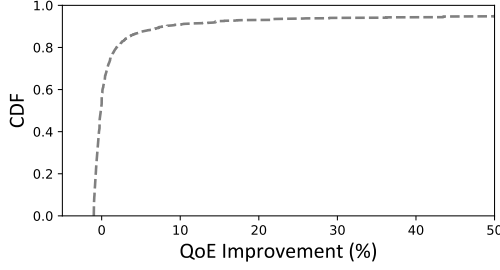


Figure 4: QoE improvements in simulation

- *high QoE with low latency*: Users like A are more sensitive to latency and experience a higher QoE peak value at low latency, indicating that investing fewer system resources in such users will result in greater QoE benefits.
- *low QoE with all latency*: Users such as B, C, and E manifest low QoE peak values, irrespective of whether they are assigned low or high latency. Thus, we can allocate fewer resources to these users, who are less sensitive to both response latency and recommendation accuracy.
- *high QoE with high latency*: Users like H show high sensitivity to recommendation accuracy and have a high QoE peak value, implying that more resources must be invested to satisfy such users.

In summary, to improve QoE while minimizing resource usage, we should prioritize resource allocation for users showing high QoE with low latency, balance resource allocation for users with high accuracy sensitivity, and assign simple strategies to users with consistently low QoE values across all latencies.

Some previous works also use user heterogeneity for resource allocation [32, 33]. The work most similar to ours is E2E [33]. However, it lacks consideration of user interests and uses a sub-optimal greedy algorithm. In contrast, our work accounts for user preferences for the personalized latency-accuracy tradeoff and proposes an algorithm that can achieve an approximate optimal solution with little additional time in online services.

2.4 QoE Improvement in simulation

We conduct a simulation experiment to prove the effectiveness of using user heterogeneity to allocate resources. We select 1350 users with recorded request data and build a simple QoE model to estimate their QoE on different strategies. Suppose the selected strategy set is S and the user set is U . We re-allocate the strategy for each user via the following steps:

1. Randomly select a user U' from U , then find the strategy S' from S that maximizes the QoE of user U' ;
2. Assign S' to user U' , and remove S' , U' from S , U respectively;
3. Select the next user from U and repeat step 1.

The results of the re-allocation demonstrate a significant improvement in the average Quality of Experience (QoE) compared to the original resource allocation. Specifically, the proposed method has increased the average QoE by 25.98%, as illustrated in Figure 4, which presents the QoE cumulative distribution function. It is worth

Symbols	Description
n	the number of concurrent requests
m	the number of strategies
I	the constraint of system overhead
B	the budget set
r_i	the i th request, or the i th user
s_j	the j th strategy
x_i	the strategy allocated to r_i
$c(\cdot)$	the cost function
$q_i(\cdot)$	the QoE model of r_i
$U(\cdot)$	the utility function

Table 2: Summary of mathematics symbols

noting that the proposed greedy allocation represents a suboptimal solution, highlighting the potential benefits of making tradeoffs between response latency and recommendation accuracy in resource allocation.

2.5 Takeaways

We summarize the observations as follows:

- We analyzed user data from four days and found a quadratic-like relationship between QoE and latency, unlike the sigmoid-like curve in traditional services. The QoE heatmap showed user heterogeneity in QoE sensitivity to response latency and recommendation accuracy.
- The simulation experiment shows that the average QoE increases by approximately 25.98% when greedily re-allocating resources according to user heterogeneity.

3 HYDRUS: MODELING

In this section, we introduce how to formulate resource allocation according to user heterogeneity as the utility maximization problem and then elaborate on how to solve this problem in polynomial time.

3.1 Formulation

Assuming that there are n concurrent requests r_1, r_2, \dots, r_n and m strategies s_1, s_2, \dots, s_m and the maximum system overhead is I . The mathematics symbols are shown in Table 2. All the allocations that satisfy the constraint of system overhead I form a budget set B :

$$B(c, I) = \{\mathbf{x} \mid \sum_{i=1}^n c(x_i) \leq I\} \quad (1)$$

where $c(\cdot)$ is the cost function used to estimate the latency under a specific strategy. Furthermore, x_i is the strategy allocated to r_i and \mathbf{x} is an n -dimensional vector, which is the strategy allocation set of n requests.

We define a *utility function* $U(\mathbf{x})$, which is a term applied in economics to model the worth or value of different available choices. In our problem, we use the overall QoE as the utility under the allocation \mathbf{x} :

$$U(\mathbf{x}) = \sum_{i=1}^n q_i(c(x_i)) \quad (2)$$

Algorithm 1 Finding the optimal strategy by solving the system of equations

Require: $q_i(\cdot)$, $c(\cdot)$, I , requests r_1, \dots, r_n

```

1: /* initialize derivative functions */
2:  $q'_1, q'_2, \dots, q'_n \leftarrow$  the derivative of  $q_i$ 
3: /* equation (6) */
4:  $f_1(\mathbf{c}) : c_1 + c_2 + \dots + c_n = I$ 
5: /* equation (7) */
6: for each  $i \in [2, \dots, n]$  do
7:    $f_i(\mathbf{c}) : q'_1(c_1) = q'_i(c_i)$ 
8: end for
9: /* solving the system of equations */
10:  $\mathbf{c}^* \leftarrow$  solving  $F(\mathbf{c}) = (f_1(\mathbf{c}), \dots, f_n(\mathbf{c}))$ 
11: /* solving the optimal strategy allocation */
12:  $\mathbf{x}^* \leftarrow$  solving  $\mathbf{c}^* = c(\mathbf{x})$ 
13: return  $\mathbf{x}^*$ 

```

where $q_i(\cdot)$ represents user r_i 's QoE model, which takes the latency of a specific strategy as input and outputs a QoE value in $[0, 1]$.

Then, the optimal strategy allocation $\mathbf{x}^*(c, I)$ is the maximum utility bundle of all bundles in the budget set:

$$\mathbf{x}^*(c, I) = \operatorname{argmax}_{\mathbf{x} \in B(c, I)} U(\mathbf{x}) \quad (3)$$

The process of finding $\mathbf{x}^*(c, I)$ is called the *utility maximization problem*. According to the description in § 2.2, $q_i(\cdot)$ is a continuous and strictly concave function, which brings the same properties to utility function $U(\cdot)$ because it is a linear combination of $q_i(\cdot)$. Therefore, $\mathbf{x}^*(c, I)$ exists and is unique since the corresponding *indifference curve* is strictly convex [17, 26].

3.2 Solution

For simplicity, we use c_i to represent the cost of user r_i under strategy x_i ; then, Equation 2 can be written as:

$$U(c_1, c_2, \dots, c_i) = \sum_i^n q_i(c_i) \quad (4)$$

Based on Walras's Law [29], we can obtain the optimal strategy allocation $\mathbf{x}^*(c, I)$ when the following conditions are met:

$$\begin{cases} \text{MU}_1 = \text{MU}_2 = \dots = \text{MU}_n \\ c_1 + c_2 + \dots + c_n = I \end{cases} \quad (5)$$

where $\text{MU}_i = \partial U / \partial c_i$ represents the *marginal utility* of user r_i at c_i . Equivalently, we can express Equation 5 as follows:

$$\frac{dq_1}{dc_1} = \frac{dq_2}{dc_2} = \dots = \frac{dq_n}{dc_n} \quad (7)$$

For a q function for which it is easy to find the derivative (e.g., the quadratic function), we can solve this system of equations through simple matrix operations in polynomial time, as shown in Algorithm 1. However, in some cases, q is a model based on deep learning or a higher-order function: a considerable amount of time is required to find the derivative of q or to solve the system of higher-order equations.

Inspired by DCAF [16], we propose Algorithm 2 to quickly produce approximately optimal solutions. It's different from DCAF that our algorithm can better solve the problem of QoE fairness, as

Algorithm 2 Finding the personalized parameter for each user to maximize the overall QoE

Require: $q_i(\cdot)$, $c(\cdot)$, \mathbf{k}_{max} , I , δ , requests r_1, \dots, r_n , strategies s_1, \dots, s_m

```

1: /* initialize global variables */
2:  $\mathbf{k}_l \leftarrow [0, \dots, 0]$ 
3:  $\mathbf{k}_r \leftarrow [\mathbf{k}_{max}^1, \dots, \mathbf{k}_{max}^n]$ 
4:  $\text{diff} \leftarrow +\infty$ 
5: /* end the loop when the overall cost is close to the limit  $I$  */
6: while  $\text{diff} > \delta$  do
7:   /* find the strategy for which the derivative is closest to  $\mathbf{k}_m^i$  */
8:   for each  $i \in [1, \dots, n]$  do
9:      $\mathbf{k}_m^i \leftarrow (\mathbf{k}_l^i + \mathbf{k}_r^i) / 2$ 
10:     $x_i \leftarrow \arg \max_j (q_i(c(s_j)) - \mathbf{k}_m^i c(s_j)), \forall j \in [1, m]$ 
11:  end for
12:   $\mathbf{x}^* \leftarrow \{x_1, \dots, x_i\}, \mathbf{c} \leftarrow \{c(s_j) | j \in \mathbf{x}^*\}$ 
13:   $\text{diff} \leftarrow |\sum(\mathbf{c}) - I|$ 
14:  if  $\sum(\mathbf{c}) < I$  then
15:    /* decrease  $k$  to improve QoE */
16:     $\mathbf{k}_r \leftarrow [\mathbf{k}_m^1, \dots, \mathbf{k}_m^n]$ 
17:  else
18:    /* increase  $k$  to reduce cost */
19:     $\mathbf{k}_l \leftarrow [\mathbf{k}_m^1, \dots, \mathbf{k}_m^n]$ 
20:  end if
21: end while
22: return  $\mathbf{k}_m$ 

```

described in § 5.2. Assuming each user's marginal utility is equal to k , $\sum_i c_i$ will increase as k decreases due to the concavity of the q function. Thus, we can increase or decrease k until $\sum_i c_i$ is sufficiently close to I using bisection search in the range of $[0, k_{max}]$. We assign a different \mathbf{k}^i to each user in the range of $[0, \mathbf{k}_{max}^i]$ rather than use a shared k , in which \mathbf{k}_{max}^i can be represented as:

$$\mathbf{k}_{max}^i \approx \max \left(\frac{q_i(c(s_1))}{c(s_1)}, \dots, \frac{q_i(c(s_m))}{c(s_m)} \right) \quad (8)$$

In large-scale systems, the number of requests per second can reach tens of thousands. Running Algorithm 2 directly to allocate the strategy for each request will introduce additional latency, which can affect the stability of online services. To avoid this, we save the solved parameter \mathbf{k}_m for calculating personalized strategy for each request in the online stage, as discussed in § 4.2.

4 HYDRUS: SYSTEM

We design Hydrus, a personalized resource allocation system, to improve the overall QoE. Hydrus allocates a recommendation strategy for each request based on user heterogeneity in QoE sensitivity to recommendation accuracy and response latency. As shown by the system architecture in Figure 5, Hydrus is composed of a parameter solver module and a strategy allocation module:

4.1 Parameter Solver

The parameter solver module subscribes to the real-time log stream published by the strategy allocation module and processes the logs within a particular time window (e.g., 15 minutes). These logs

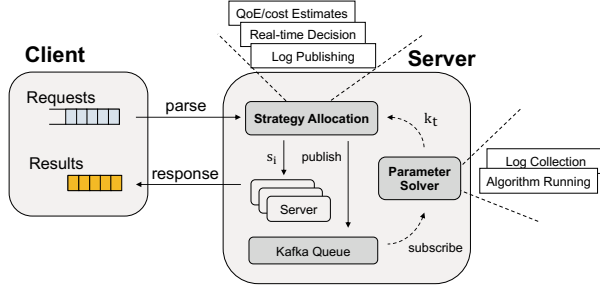


Figure 5: Hydrus architecture

contain all the information described in § 4.2. First, the module adds the system resource allocated to each request to obtain the amount of system allocatable resources I , which represents the maximum load capacity of the system in this time window. Notably, system resource is an abstract concept that can represent any resource cost target, such as CPU time, memory percentage, and so forth. The specific target depends on the business needs and concerns. Then, the module calculates the parameter \mathbf{k} that can maximize user QoE in this time window according to Algorithm 2. Finally, parameter \mathbf{k} can be saved for the next online decision to allocate an optimal recommendation strategy for each request.

Time-segmented Parameters. In a large-scale online serving systems, queries per second (QPS) is often used to represent the amount of system traffic received in one second. Generally, QPS changes over time throughout the day and reaches a peak at night, resulting in the sum of resource consumption varying accordingly. Therefore, the parameter \mathbf{k}_t in every time window is also different due to the difference in the system resource limit I_t . We take 15 minutes as the time window and calculate the corresponding parameter \mathbf{k}_t . Therefore, the set of all parameters in a day n is:

$$S^n = \{\mathbf{k}_t^n | 0 \leq t \leq 95\} \quad (9)$$

We update the parameters every 15 minutes and use \mathbf{k}_t to calculate the optimal strategy for requests coming in the next time window based on the assumption that the system resource limits of adjacent time windows are the same. Maintaining a high update frequency guarantees the maximum probability of allocating the optimal strategy for users. Alternatively, we could also update the parameters once a day to save computing resources at the cost of decreased accuracy. For example, \mathbf{k}_t^n can be used to determine the allocation strategy for the requests within the time window t of day $n + 1$. Which update strategy is used depends on the budget for computing resources.

4.2 Strategy Allocation

The online decision module first estimates the request's QoE value and resource cost for each strategy. Then, it uses parameter \mathbf{k}_t to allocate the optimal strategy for each request to maximize the overall QoE in the current time window. Finally, it publishes streams of request logs using Kafka[18], including all the data required by the parameter solver.

Algorithm 3 Finding the optimal strategy for each user in a time window t

Require: $q_i(\cdot)$, $c(\cdot)$, \mathbf{k}_t , requests r_1, \dots, r_n , strategies s_1, \dots, s_m

- 1: */* obtain the parameter \mathbf{k} of time window t */*
- 2: $\mathbf{k} \leftarrow \mathbf{k}_t$
- 3: */* find the optimal strategy with \mathbf{k} */*
- 4: **for** each $i \in [1, \dots, n]$ **do**
- 5: $p_{max} = 0$
- 6: **for** each $j \in [1, \dots, m]$ **do**
- 7: $p = q_i(c(s_j)) - \mathbf{k}^T c(s_j)$
- 8: **if** $p > p_{max}$ **then**
- 9: $p_{max} = p$
- 10: $s_{max} = s_j$
- 11: **end if**
- 12: **end for**
- 13: $x_i = s_{max}$
- 14: **end for**
- 15: $\mathbf{x}^* \leftarrow \{x_1, \dots, x_i\}$
- 16: **return** \mathbf{x}^*

QoE and cost estimates. We use VCR to estimate user QoE in § 2.1. Besides VCR, we could use many other targets to quantify user QoE in a large-scale recommendation system, such as watch time and click-through rate. Like the system resources mentioned above, the selection of a QoE target depends on the specific scenario and requirements. As the user QoE value, we choose the output of the existing online deep-learning model that has been trained with a large amount of user data and provides accurate predictions. However, the estimated QoE values of online models are produced under the same default strategy (e.g., strategy s_a), so we cannot obtain QoE directly for other strategies (e.g., strategies s_b , s_c , and s_d). For convenience, we use A/B testing to obtain the proportional coefficient of the other strategies relative to the default strategy. Specifically, we divide the users into four groups and use strategies s_a , s_b , s_c , and s_d to recommend videos. After running the experiment for a week (the greater the number of users is, the higher the confidence), we summarize the results and calculate the coefficient. For example, the average watch time under strategy s_a is 20 s, and that under strategies s_b , s_c , and s_d is 22 s, 24 s, and 20 s. Therefore, we can calculate the proportional coefficients to be 1.1, 1.2, and 1.0 relative to strategy s_a . Finally, we multiply these coefficients and the outputs of the online model to estimate the QoE value of different strategies. For resource cost estimates, we also use A/B testing to summarize the proportional coefficient and calculate the resource cost of different strategies. Unlike user QoE estimates, resource cost is related to a specific strategy but not related to users.

Real-time Decision. The first step is to obtain the matching parameter \mathbf{k}_t from storage based on the current timestamp. Then, we can then estimate the request's QoE value and resource cost for every strategy using the above-described approaches. According to Algorithm 3, the optimal strategy s_{max} is allocated to each request. Step 7 is the essential operation that selects the strategy with the best marginal benefit for request r_i with parameter \mathbf{k}^i .

Log Publishing. After allocating the optimal strategy for each request, we publish all the pertinent logs to a particular Kafka message queue. This message queue is subscribed to by the parameter

Listing 1 The data structure of log information

```

struct ValueInfo {
    int strategy;
    double cost;
    double qoe;
};

struct LogInfo {
    int request_id;
    int user_id;
    int device_id;
    vector<ValueInfo> value_info;
};

```

solver, which will then utilize these logs for the subsequent parameter computation. The log's data structure is composed of *ValueInfo* and *LogInfo*, as shown in Listing 1.

ValueInfo is a structure that stores the estimated QoE value and resource cost associated with a particular recommendation strategy. A vector is used to hold *ValueInfo* for various strategies in *LogInfo*.

5 EVALUATION

In this section, we evaluate the effectiveness of our method in terms of QoE gain, latency reduction, throughput, and QoE fairness in trace-driven simulation. Then, we conduct online experiments to validate Hydrus's accuracy, efficiency, and tolerance in Kuaishou's large-scale recommendation system.

5.1 Experimental setup

Data. The QoE model requires a large amount of training data for each user. To improve the accuracy and confidence of the experiments, we clean and filter hundreds of millions of trace data in the dataset. However, most users in the dataset have only a small amount of request data that will cover only some strategies. Therefore, we filter users with fewer than 100 instances of recorded trace data. Finally, we obtain a training dataset that contains approximately 0.3M users and 42M trace data.

Methods. We illustrate the advantages of Hydrus by comparing it with the following methods:

- **Default strategy** refers to the original strategy used in Kuaishou's online short-form video service.
- **Greedy strategy** is the strategy used in § 2.4 that prioritizes users who arrive earlier to be allocated more resources.
- **DCAF** [16] is a dynamic computation allocation framework that allocates computing resources according to the importance of users.

Metric. We use QoE gain, latency reduction, throughput, and QoE fairness to demonstrate the effectiveness of Hydrus, as described below:

- **QoE gain** is the relative improvement in average QoE over that of the default strategy, which is defined as $(Q_{\text{new}} - Q_{\text{ori}})/Q_{\text{ori}}$.
- **Latency reduction** is the relative reduction in overall latency compared to that of the default strategy, which is defined as $(L_{\text{ori}} - L_{\text{new}})/L_{\text{ori}}$.
- **Throughput** is the maximum request number the system can handle given the limitation of overall system resources.

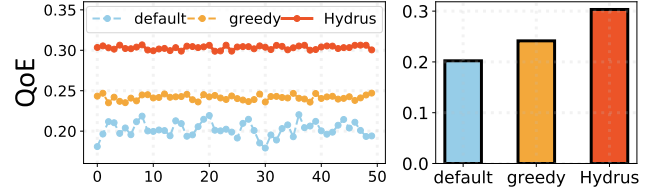


Figure 6: The average QoE obtained from 50 experiments using different methods. The average QoE of Hydrus is 49% higher than that of the default strategy and 25% higher than that of the greedy strategy.

- **QoE Fairness** is used to quantify the discrimination of users by considering the standard deviation σ of the latency difference.

5.2 Trace-driven experiment

As described in § 3, we first fit a function $c(\cdot)$ to obtain the latency of different strategies and then build a QoE model $q_i(\cdot)$ for each user r_i . Then, we create a latency set \mathcal{L} and QoE matrix \mathbf{Q} to facilitate the calculation of Algorithms 2 and 3:

- \mathcal{L} : We use m to represent the number of unique strategies in the dataset, so the mathematical definition of \mathcal{L} is given in Equation 10. For convenience, we use \mathcal{L}_j to represent the latency of strategy s_j .

$$\mathcal{L} = \{c(s_j) | j \in [1, m]\} \quad (10)$$

- \mathbf{Q} : The definition of \mathbf{Q} is given by Equation 11, in which n refers to the user number in an experiment. Therefore, Q_{ij} represents the estimated QoE value for strategy s_j and user r_i .

$$\mathbf{Q} = \begin{pmatrix} Q_1(\mathcal{L}_1) & \cdots & Q_1(\mathcal{L}_k) \\ \vdots & \ddots & \vdots \\ Q_n(\mathcal{L}_1) & \cdots & Q_n(\mathcal{L}_k) \end{pmatrix} \quad (11)$$

QoE gain. We run 50 experiments to illustrate that Hydrus outperforms the default and greedy strategy. In each experiment, we randomly select 1000 users and their trace data from the dataset. We use Q_{ori} to represent the average QoE of the original strategy allocation. Then, we apply the greedy strategy and the improved algorithm in Hydrus to re-allocate the strategies: the QoE values are expressed as Q_{greedy} and Q_{Hydrus} , respectively. The experimental results are illustrated in Figure 6. From the line chart, we can observe that Q_{Hydrus} is stably maintained at approximately 0.3, which is higher than Q_{ori} and Q_{greedy} in these 50 experiments. Furthermore, we calculate the average Q_{ori} , Q_{greedy} and Q_{Hydrus} in these 50 experiments, as shown in the bar chart: Hydrus increases QoE by 49% and 25% compared with that of the default and greedy strategy, which illustrates the effectiveness of Hydrus.

The tradeoff between QoE and latency. In this experiment, we explore the relationship between QoE and latency by dynamically adjusting k when we make a tradeoff between them. It is worth noting that k represents the average value of \mathbf{k}_m used in Algorithm 2. Specifically, we choose a subset containing 1000 users as the baseline, of which the average QoE is 0.21 and the average latency is 435 ms. We control k to increase from 0 with a 1×10^{-3} interval until k reaches the maximum value and solves the optimal allocation

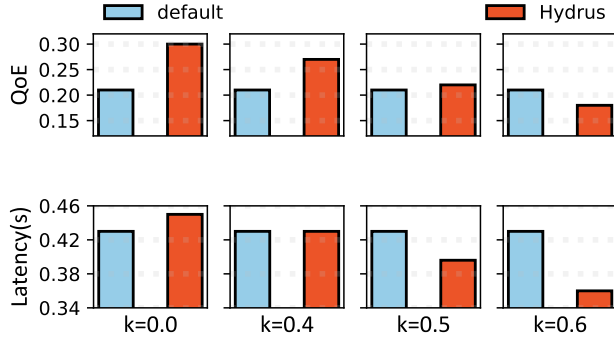


Figure 7: The change in average QoE and average latency with the increase in k . The magnitude of k is 10^{-3} .

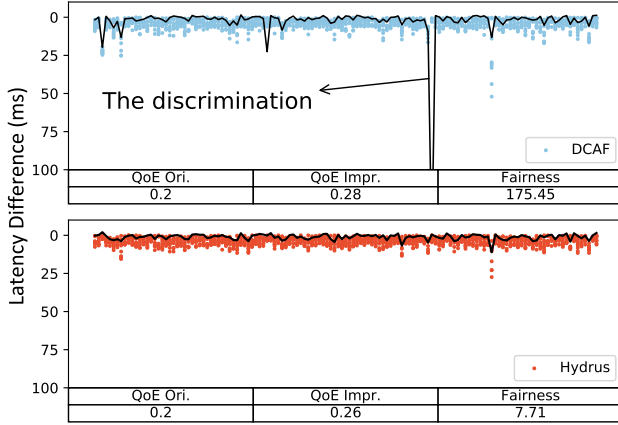


Figure 8: The difference in latency compared to the default strategy using the allocation by DCAF and Hydrus from 10 experiments for every user.

for each k . We select 4 most representative experimental results when k is 0, 0.4×10^{-3} , 0.5×10^{-3} and 0.6×10^{-3} , as shown in Figure 7. $k = 0$ indicates assigning each user their favorite strategy; however, the average latency in this case exceeds the default latency limit. As k increases to 0.4×10^{-3} , Hydrus improves the QoE by approximately 35.6%, and the average latency is equal to the limit. We continue to increase k until the average QoE of Hydrus is almost the same as that of the default strategy; in this case, the average latency is reduced by approximately 10.1%. Finally, the QoE value and latency of Hydrus are lower than that of the default strategy.

QoE Fairness. Providing QoE fairness is critical because user QoE will decrease if the latency difference between two adjacent requests is perceivable. In this experiment, we randomly select 10 historical requests of each user to run 10 experiments. In Figure 8, the horizontal axis represents different users, and the vertical axis represents the latency difference $|L_{DCAF} - L_{ori}|$ and $|L_{Hydrus} - L_{ori}|$. The black line chart records the average latency difference for the 10 experimental results of each user. DCAF sacrifices the QoE of

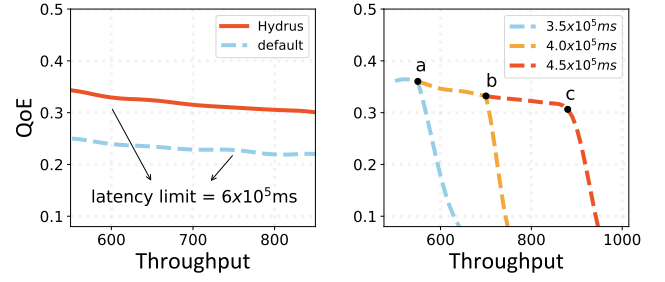


Figure 9: The QoE value under different loads.

some users to obtain the global maximum QoE. By contrast, Hydrus pursues the improvement of QoE and fairness among users. The variance of the latency difference decreases from 119.24 to 7.99 at the cost of a small decrease in QoE improvement, which indicates that Hydrus can reduce the discrimination significantly.

Throughput. Figure 9 shows the advantage of Hydrus in improving throughput compared with the default strategy or under different latency limits. First, we evaluate the throughput performance of Hydrus and the default strategy under the limit of 6×10^5 ms, which is the maximum latency limit in the dataset. Hydrus produces an allocation with the maximum user QoE in all the strategies. As shown in the left figure, we increased the request load from 550 to 850 to simulate traffic changes over time in one day. The QoE values of Hydrus and the default strategy decrease from 0.34 to 0.30 and 0.25 to 0.22, which shows that Hydrus can achieve at least 54.5% higher throughput than the default strategy without a decrease in QoE. Then, we demonstrate the automatic downgrade mechanism of Hydrus in the right figure when the number of requests increases. When the overall latency exceeds the limit, Hydrus downgrades user QoE to reduce the failure rate and ensure the stability of the service (the inflection points a , b , and c). In addition, we can actively control the latency limit according to the number of requests so that Hydrus can adaptively adjust service quality to achieve a tradeoff between user QoE and system resources.

5.3 Online evaluation in a production system

We evaluate Hydrus's capability to improve QoE and reduce resource costs in three products of Kuaishou and two recommendation scenarios. The result is presented in Table 3 and 4, and more details about the experimental setup are described next.

QoE improving. We apply Hydrus to live streaming recommendation services in two products of Kuaishou: Kuaishou Flagship² and Kuaishou Express³. We select 20% of users to use Hydrus for recommendation strategy allocation, and the remaining 80% of users use the default recommendation strategy as the baseline. The experimental and control groups use the same amount of system resources in the same time window for a week. The results show that CTR and watching time increase by 1.034% and 1.753% in Kuaishou Flagship. Furthermore, CTR increases by 1.362% and watching time increases by 1.851% in Kuaishou Express. These results demonstrate that Hydrus achieves a substantial improvement in user QoE.

²A global leading short video and content-based social networking platform

³A variant of Kuaishou Flagship with a simpler interface

	CTR	Watching Time
Kuaishou Flagship	+1.03%	+1.75%
Kuaishou Express	+1.36%	+1.85%

Table 3: Improvement in QoE in livecast recommendation.

	Average Latency	CPU utilization
Kuaishou Express	-3.04%	-4.01%
Kwai	-5.71%	-7.71%

Table 4: Resource saving in video recommendation.

Resources saving. In § 5.2, we conduct offline experiments showing that the latency decreases as parameter k increases, which means fewer system resources are required. When the overall QoE is the same, Hydrus requires far fewer resources than does the default allocation strategy, as shown in Figure 7. To prove this finding holds in the online system, we conduct a 7-day experiment on 20% of users in the video recommendation system of Kuaishou Express and Kwai⁴. The results show that there are benefits in average latency and CPU utilization. Hydrus reduces the average latency by approximately 3% to 4% and lowers CPU utilization by approximately 5% to 7%.

Light-weight system. Since most of the time-consuming calculations are in the parameter solver module, the online module of Hydrus introduces little additional running time. As shown in Table 5, the average running time is only approximately 5 ms, which will not bring any system pressure.

6 RELATED WORK

In this section, we survey and present work related to improving user QoE and making tradeoffs between latency and quality in Internet services. Additionally, we demonstrate the differences and advantages of Hydrus over existing works.

QoE improvement. How to improve user QoE in various Internet services, such as web [2, 10, 33], livecast streaming [32], video streaming [3, 11], and mobile applications [8, 24], has been studied for many years. As described in prior work [1, 5, 9, 13, 23, 24, 33], QoE usually decreases as a sigmoid-like curve as latency increases. Therefore, most work aims to improve QoE by cutting tail/average latency.

Some works use a common approach to cut tail latency by sending redundant requests and using the first completed request [27, 28, 30, 31]. When requests cannot be replicated, some scheduling techniques are proposed to cut long tail latencies (e.g., prioritization strategies [35], load balancing techniques [15]). In contrast to these works, RobinHood [4] utilizes the cache to reduce overall request tail latency. C3 [25] prefers faster replicas along with distributed rate control and backpressure to reduce tail latency.

Moreover, some methods are proposed to improve QoE based on user heterogeneity and preferences. KLOTSKI [6] prioritizes the web content most relevant to a user's preferences to improve

	Running Time(ms)
Kuaishou Flagship	~5.31
Kuaishou Express	~5.67
Kwai	~4.52

Table 5: Additional running time introduced in the online system.

user experience. E2E [33] employs a greedy algorithm to prioritize resource allocation to users with high QoE sensitivity. Similarly, HeteroCast [32] schedules CDNs based on the QoE sensitivity to QoS metrics in large-scale livecast scenarios. Different from the above methods, DCAF [16] formulates the resource allocation as a knapsack problem to maximize the overall user value.

Unlike the sigmoid-like relationship between latency and QoE described in prior work, Hydrus addresses the novel quadratic-like relationship in short-form video services by making tradeoffs between response latency and recommendation accuracy.

Latency-quality tradeoffs. With a latency deadline, Zeta [14] allocates processor time among service requests to maximize the quality and minimize the variance of the response. It can free resources for those requests that might have timed out and produced no results. Timecard [24] tracks delays across activities and estimates network transfer times to adapt its processing time and control the end-to-end delay for requests. In addition to these deadline-driven methods, DQBarge [9] is proactive in enabling better tradeoffs by propagating critical information along the causal path of request processing.

However, user preferences are excluded from consideration in these methods. Hydrus can achieve optimal resource allocation to users by making use of their different QoE sensitivity to response latency and recommendation accuracy, resulting in better overall QoE and lower latency.

7 CONCLUSION

We find a novel quadratic-like relationship between user QoE and response latency, as well as user heterogeneity in QoE sensitivity to response latency and recommendation accuracy. Based on these findings, we propose Hydrus, a resource allocation system that aimed to improve personalized QoE by making tradeoffs between response latency and recommendation accuracy. Our experiments on Kuaishou, a popular video app with hundreds of millions of users, demonstrate that Hydrus can significantly increase QoE or reduce latency without sacrificing QoE. Overall, our research shows the effectiveness of Hydrus in addressing the challenges faced by short-form video services and highlights the importance of personalized resource allocation in enhancing the user experience.

ACKNOWLEDGMENTS

We express our gratitude to Feng Gao, Jiayuan Xing, Lijun Liu for their valuable contributions to the development of Hydrus. We also thank Feng Jiang, Lifei Zhu, Sukun Zhang for their insightful discussions and helpful advice.

⁴An international version of Kuaishou Flagship

REFERENCES

- [1] Ioannis Arapakis, Xiao Bai, and B Barla Cambazoglu. 2014. Impact of response latency on user behavior in web search. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. 103–112.
- [2] Athula Balachandran, Vaneet Aggarwal, Emir Halepovic, Jeffrey Pang, Srinivasan Seshan, Shobha Venkataraman, and He Yan. 2014. Modeling web quality-of-experience on cellular networks. In *Proceedings of the 20th annual international conference on Mobile computing and networking*. 213–224.
- [3] Athula Balachandran, Vyas Sekar, Aditya Akella, Srinivasan Seshan, Ion Stoica, and Hui Zhang. 2013. Developing a predictive model of quality of experience for internet video. *ACM SIGCOMM Computer Communication Review* 43, 4 (2013), 339–350.
- [4] Daniel S Berger, Benjamin Berg, Timothy Zhu, Siddhartha Sen, and Mor Harchol-Balter. 2018. Robinhood: Tail latency aware caching–dynamic reallocation from cache-rich to cache-poor. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*. 195–212.
- [5] Enrico Bocchi, Luca De Cicco, and Dario Rossi. 2016. Measuring the quality of experience of web users. *ACM SIGCOMM Computer Communication Review* 46, 4 (2016), 8–13.
- [6] Michael Butkiewicz, Daimeng Wang, Zhe Wu, Harsha V Madhyastha, and Vyas Sekar. 2015. Klotki: Reprioritizing web content to improve user experience on mobile devices. In *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*. 439–453.
- [7] Craig Chambers, Ashish Raniwala, Frances Perry, Stephen Adams, Robert R Henry, Robert Bradshaw, and Nathan Weizenbaum. 2010. FlumeJava: easy, efficient data-parallel pipelines. *ACM Sigplan Notices* 45, 6 (2010), 363–375.
- [8] Qi Alfred Chen, Haokun Luo, Sanae Rosen, Z Morley Mao, Karthik Iyer, Jie Hui, Kranthi Sontineni, and Kevin Lau. 2014. Qoe doctor: Diagnosing mobile app qoe with automated ui control and cross-layer analysis. In *Proceedings of the 2014 Conference on Internet Measurement Conference*. 151–164.
- [9] Michael Chow, Mosharaf Chowdhury, Kaushik Veeraraghavan, Christian Cachin, Michael Cafarella, Wonho Kim, Jason Flinn, Marko Vukolić, Sonia Margulis, Inigo Goiri, et al. 2016. Dqbarge: Improving data-quality tradeoffs in large-scale internet services. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. 771–786.
- [10] Diego Neves da Hora, Alemnew Sheferaw Asrese, Vassilis Christophides, Renata Teixeira, and Dario Rossi. 2018. Narrowing the gap between QoS metrics and Web QoE using Above-the-fold metrics. In *International Conference on Passive and Active Network Measurement*. Springer, 31–43.
- [11] Florin Dobrian, Vyas Sekar, Asad Awan, Ion Stoica, Dilip Joseph, Aditya Ganjam, Jibin Zhan, and Hui Zhang. 2011. Understanding the impact of video quality on user engagement. *ACM SIGCOMM computer communication review* 41, 4 (2011), 362–373.
- [12] Fahad R Dogar, Thomas Karagiannis, Hitesh Ballani, and Antony Rowstron. 2014. Decentralized task-aware scheduling for data center networks. *ACM SIGCOMM Computer Communication Review* 44, 4 (2014), 431–442.
- [13] Qingzhu Gao, Prasenjit Dey, and Parvez Ahammad. 2017. Perceived performance of top retail webpages in the wild: Insights from large-scale crowdsourcing of above-the-fold qoe. In *Proceedings of the Workshop on QoE-based Analysis and Management of Data Communication Networks*. 13–18.
- [14] Yuxiong He, Sameh Elnikety, James Larus, and Chenyu Yan. 2012. Zeta: Scheduling interactive services with partial execution. In *Proceedings of the Third ACM Symposium on Cloud Computing*. 1–14.
- [15] Seyyed Ahmad Javadi and Anshul Gandhi. 2017. Dial: Reducing tail latencies for cloud applications via dynamic interference-aware load balancing. In *2017 IEEE International Conference on Autonomic Computing (ICAC)*. IEEE, 135–144.
- [16] Biye Jiang, Pengye Zhang, Rihan Chen, Xinchun Luo, Yin Yang, Guan Wang, Guorui Zhou, Xiaoqiang Zhu, and Kun Gai. 2020. DCAF: A Dynamic Computation Allocation Framework for Online Serving System. *arXiv preprint arXiv:2006.09684* (2020).
- [17] David M Kreps. 2013. *Microeconomic foundations I: choice and competitive markets*. Vol. 1. Princeton university press.
- [18] Jay Kreps, Neha Narkhede, Jun Rao, et al. 2011. Kafka: A distributed messaging system for log processing. In *Proceedings of the NetDB*, Vol. 11. 1–7.
- [19] Kuaishou. 2023. Kuaishou Technology Announces Fourth Quarter and Full Year 2022 Financial Results. <https://ir.kuaishou.com/news-releases/news-release-details/kuaishou-technology-announces-fourth-quarter-and-full-year-2022>. Accessed: 2023-03-29.
- [20] Conglong Li, David G Andersen, Qiang Fu, Sameh Elnikety, and Yuxiong He. 2017. Workload analysis and caching strategies for search advertising systems. In *Proceedings of the 2017 Symposium on Cloud Computing*. 170–180.
- [21] Conglong Li, David G Andersen, Qiang Fu, Sameh Elnikety, and Yuxiong He. 2018. Better caching in search advertising systems with rapid refresh predictions. In *Proceedings of the 2018 World Wide Web Conference*. 1875–1884.
- [22] Hongzi Mao, Malte Schwarzkopf, Shaileshh Bojja Venkatakrishnan, Zili Meng, and Mohammad Alizadeh. 2019. Learning scheduling algorithms for data processing clusters. In *Proceedings of the ACM Special Interest Group on Data Communication*. 270–288.
- [23] Ravi Netravali, Vikram Nathan, James Mickens, and Hari Balakrishnan. 2018. Vesper: Measuring time-to-interactivity for web pages. In *15th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 18)*. 217–231.
- [24] Lenin Ravindranath, Jitendra Padhye, Ratul Mahajan, and Hari Balakrishnan. 2013. Timecard: Controlling user-perceived delays in server-based mobile applications. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*. 85–100.
- [25] Lalith Suresh, Marco Canini, Stefan Schmid, and Anja Feldmann. 2015. C3: Cutting tail latency in cloud data stores via adaptive replica selection. In *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*. 513–527.
- [26] Hal R Varian and Hal R Varian. 1992. *Microeconomic analysis*. Vol. 3. Norton New York.
- [27] Ashish Vulimiri, Philip Brighten Godfrey, Radhika Mittal, Justine Sherry, Sylvia Ratnasamy, and Scott Shenker. 2013. Low latency via redundancy. In *Proceedings of the ninth ACM conference on Emerging networking experiments and technologies*. 283–294.
- [28] Ashish Vulimiri, Oliver Michel, P Brighten Godfrey, and Scott Shenker. 2012. More is less: reducing latency via redundancy. In *Proceedings of the 11th ACM Workshop on Hot Topics in Networks*. 13–18.
- [29] Leon Walras. 2013. *Elements of pure economics*. Routledge.
- [30] Zhe Wu, Curtis Yu, and Harsha V Madhyastha. 2015. Costlo: Cost-effective redundancy for lower latency variance on cloud storage services. In *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*. 543–557.
- [31] Neeraja J Yadwadkar, Ganesh Ananthanarayanan, and Randy Katz. 2014. Wrangler: Predictable and faster jobs using fewer resources. In *Proceedings of the ACM Symposium on Cloud Computing*. 1–14.
- [32] Rui-Xiao Zhang, Ming Ma, Tianchi Huang, Hanyu Li, Jiangchuan Liu, and Lifeng Sun. 2020. Leveraging QoE Heterogeneity for Large-Scale Livechat Scheduling. In *Proceedings of the 28th ACM International Conference on Multimedia*. 3678–3686.
- [33] Xu Zhang, Siddhartha Sen, Danian Kurniawan, Haryadi Gunawi, and Junchen Jiang. 2019. E2E: embracing user heterogeneity to improve quality of experience on the web. In *Proceedings of the ACM Special Interest Group on Data Communication*. 289–302.
- [34] Hao Zhou, Ming Chen, Qian Lin, Yong Wang, Xiaobin She, Sifan Liu, Rui Gu, Beng Chin Ooi, and Junfeng Yang. 2018. Overload control for scaling wechat microservices. In *Proceedings of the ACM Symposium on Cloud Computing*. 149–161.
- [35] Timothy Zhu, Alexey Tumanov, Michael A Kozuch, Mor Harchol-Balter, and Gregory R Ganger. 2014. Prioritymeister: Tail latency qos for shared networked storage. In *Proceedings of the ACM Symposium on Cloud Computing*. 1–14.