

# LinkedOut: Linking World Knowledge Representation Out of Video LLM for Next-Generation Video Recommendation

## Supplementary Material

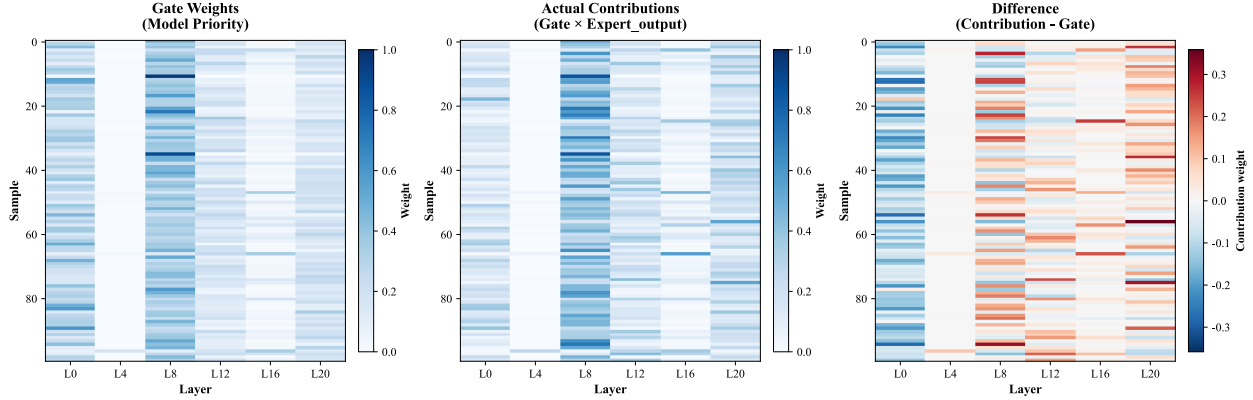


Figure 1: Sample-level heatmap of layer-wise MoE gate contributions. Left: **Gate Weights** ( $\pi_\ell$ ) show learned layer priorities. Middle: **Actual Contributions** ( $\pi_\ell \times \mathbf{h}^{(\ell)}$ ) reveal effective semantic utilization. Right: **Difference** highlights priority-impact discrepancy. Each row is a video sample; columns are transformer layers (L0, L4, L8, L12, L16, L20). L8 exhibits dominant contributions while L4 and L16 show minimal activation, confirming mid-level features balance visual details and abstract concepts better than the deepest layers alone.

## 6. Dataset Details

We evaluate LinkedOut on MicroLens-50K and MicroLens-100K [27], the only publicly available micro-video datasets providing raw video files rather than pre-extracted features, making them ideal for evaluating foundation-model-driven approaches. MicroLens-50K contains 50,000 users, 19,220 videos, and 359,708 interactions, while MicroLens-100K scales to 100,000 users, 19,738 videos, and 719,405 interactions. Both datasets span 15,580 fine-grained hashtag-based category tags (e.g., “#vegan\_cooking”, “#urban\_photography”) covering diverse domains including food, sports, travel, entertainment, education, and lifestyle, which align well with video foundation models’ world knowledge.

The datasets exhibit several key characteristics that make them suitable for evaluating video recommendation systems. Each video includes rich multimodal content: full-length video frames (average duration 161 seconds, with most videos under 400 seconds), audio tracks capturing speech and background music, textual titles with embedded category tags, cover images, and engagement metadata (likes, views, timestamps). This multimodal richness enables comprehensive content understanding beyond visual appearance alone. User-item interactions are timestamped comment behaviors that naturally encode sequential viewing patterns, with most users having 5-15 interactions cap-

turing realistic short-term preference dynamics. For sequential models, we limit the maximum sequence length to 13, which covers approximately 95% of user behaviors. Item popularity follows a long-tail distribution typical of real-world recommender systems, and the datasets naturally exhibit cold-start scenarios for new users and items, making them challenging testbeds where collaborative filtering signals are often sparse.

The data was collected over one year (June 2022-2023) from a public micro-video platform using a three-stage process. First, seed videos from the homepage with more than 10,000 likes were collected to ensure quality and diversity. Second, the video set was expanded by randomly sampling 10 related videos from each video’s webpage. Third, videos were filtered based on quality criteria: text length  $\geq 3$  characters (after removing meaningless symbols), image cover with  $< 75\%$  single-color area, and video file size  $> 100\text{KB}$ . User-item interactions are based on public comment data (approximately 1 comment per 100 likes), and all user and video IDs have been anonymized to protect privacy. We follow the leave-one-out evaluation protocol: interactions are sorted by timestamp, with the last interaction per user reserved for testing, the second-to-last for validation, and all earlier interactions for training. This temporal splitting reflects realistic recommendation scenarios where models predict future interactions based on historical behavior.

## 7. Implementation Details

To facilitate reproducibility, we provide comprehensive hyperparameter configurations in Table 6. Our implementation consists of three main components: the VLLM backbone for feature extraction, the Cross-Layer Knowledge-Fusion MoE for multi-level semantic aggregation, and the recommendation module for ranking. The VLLM backbone (LLaVA-OneVision 7B) remains frozen during training, and we extract intermediate representations at every 4 transformer layers to balance computational cost and semantic coverage. The MoE module learns adaptive gating weights to fuse 6 layer-wise experts, with each expert compressing token sequences to a fixed budget of 64 tokens for efficiency. Training is conducted on 8 NVIDIA H100 GPUs with mixed precision (fp16) for 100 epochs, using parameter-group-specific learning rates to stabilize convergence. All experiments use fixed random seeds (seed=42) to ensure reproducibility.

## 8. Layer-Wise Contribution Analysis

Figure 1 provides a sample-level visualization of how the Cross-Layer Knowledge-Fusion MoE distributes learned weights across transformer layers during inference. The heatmap comprises three complementary views. **Gate Weights** (left panel) visualize the soft-gating distribution  $\pi_\ell$  learned by the MoE gating network for each validation sample, reflecting the model’s learned prior on which layers contain task-relevant semantics. **Actual Contributions** (middle panel) show the weighted expert outputs  $\pi_\ell \times \mathbf{h}^{(\ell)}$ , capturing the effective semantic signal after applying gate weights to layer-wise features. **Difference** (right panel) reveals the discrepancy between assigned priority and actual contribution magnitude, exposing cases where high gate weights yield modest outputs or vice versa. This three-way decomposition enables fine-grained analysis of how the MoE module learns to prioritize and utilize different semantic levels, providing insights beyond aggregate statistics.

Across diverse video samples, we observe consistent patterns: L8 (middle layers) dominates contributions with high activation intensities (darker blue regions), while L4 and L16 remain largely inactive (lighter regions). This sample-level evidence corroborates the aggregate statistics presented in the main paper, where L8 contributes 40.9% on average, L0 16.4%, L20 20.5%, and L4 and L16 less than 7% combined. The heatmap also reveals variation across samples: certain videos activate L0 more strongly (likely those requiring fine-grained visual details), while abstract or high-level content emphasizes L20. Notably, the discrepancy between gate weights and actual contributions is most pronounced for L4 and L16, suggesting that these layers may encode information that requires larger weights to manifest effectively in the final representation.

Table 6. Complete hyperparameter configurations for LinkedOut. Parameters are grouped by component: VLLM feature extraction, MoE architecture, training procedure, and data preprocessing.

Hyperparameter	Value
<b>VLLM Backbone</b>	
Model architecture	LLaVA-OneVision 7B
Parameters frozen	Yes
Layer extraction interval	Every 4 layers (L0, L4, L8, L12, L16, L20)
Total extracted layers	6
Image resolution	$224 \times 224$
<b>MoE Architecture</b>	
Number of experts	6 (one per layer)
Token compression budget	64 tokens per layer
Gating mechanism	Soft gating (learnable weights)
MoE fusion dimension	4096 (LLaVA hidden size)
<b>Training Configuration</b>	
GPUs	$8 \times$ NVIDIA H100 (80GB)
Precision	Mixed (fp16)
Batch size	256 sequences
Training epochs	100
Optimizer	AdamW
Learning rate (video enc.)	0.0001
Learning rate (text/img enc.)	0.0001
Learning rate (rec. head)	0.00001
Weight decay	0.1
Gradient clipping norm	5.0
Learning rate scheduler	step_schedule _with_warmup
Warmup steps	0
Scheduler gap	1 epoch
Scheduler alpha	1.0
<b>Data Preprocessing</b>	
Max user history length	10 videos
Max sequence length	13 interactions
Evaluation protocol	Leave-one-out
Random seed	42

This adaptive layer selection validates our design hypothesis that video recommendation benefits from selectively fusing multi-level semantic features rather than relying solely on the final layer of a VLLM, as is common practice in vision-language models.

## 9. Efficiency Analysis

### 9.1. Feature Storage Requirements

LinkedOut’s Store-and-Retrieve architecture requires pre-computing and storing video features. Each video’s representation consists of 6 layers (L0, L4, L8, L12, L16, L20),

Table 7. Feature storage requirements (float16 precision, 6 layers, 64 tokens per layer, dimension 4096).

Dataset Scale	Videos	Storage
MicroLens-50K [27]	19,220	60.5 GB
MicroLens-100K [27]	19,738	62.2 GB
Extrapolated (1M)	1,000,000	3.15 TB
Extrapolated (10M)	10,000,000	31.5 TB

Table 8. Estimated serving capacity based on 5.964 ms online latency per query.

GPUs	Queries/Second	Concurrent Users (est.)
1× H100	168	~10K
4× H100	672	~40K
8× H100	1,344	~80K
32× H100	5,376	~320K

each with 64 compressed tokens of dimension 4096. Using float16 precision, the storage footprint per video is calculated as:

$$\text{Storage per video} = 6 \times 64 \times 4096 \times 2 \text{ bytes} = 3.15 \text{ MB} \quad (4)$$

For MicroLens-100K (19,738 videos), total storage is approximately 62.2 GB, which fits comfortably in modern server RAM or SSDs. Table 7 shows storage requirements across different scales, demonstrating that LinkedOut remains practical even for large-scale deployments with millions of videos.

**Storage optimization strategies.** The footprint can be further reduced through: (1) float16  $\rightarrow$  int8 quantization (50% reduction); (2) reducing token budget from 64 to 32 (50% reduction); (3) selecting fewer layers (e.g., only L0, L8, L20, reducing to 1.57 MB per video).

## 9.2. Real-Time Serving Capacity

Based on the measured online inference latency of 5.964 ms per user query (reported in the main paper), LinkedOut can theoretically serve  $1000/5.964 \approx 168$  queries per second on a single GPU. This translates to supporting approximately 10,000 concurrent users with a 60ms response time budget (assuming 10% overhead for batching and network latency). Table 8 shows that with horizontal scaling across multiple GPUs, LinkedOut can easily handle production-scale workloads serving millions of daily active users.

## 9.3. Model Complexity Comparison

Table 9 compares trainable parameters across methods. LinkedOut freezes the 7B-parameter VLLM backbone, training only the lightweight MoE module (approximately 6.3M parameters: 6 experts with token compressor and gating network) and recommendation head (2.1M parameters),

Table 9. Model complexity comparison. LinkedOut trains significantly fewer parameters by freezing the VLLM backbone.

Method	Trainable Params	Total Params
SASRec (ID-only)	2.1 M	2.1 M
GRU4RecV	45 M	45 M
MMGCN <sub>ID+V</sub>	38 M	38 M
LinkedOut (ours)	<b>8.4 M</b> (MoE + head)	<b>8.4 M</b> + 7B (frozen & 1-time)

totaling 8.4M trainable parameters. This is 5-10 $\times$  fewer than end-to-end VideoRec approaches that fine-tune entire video encoders (typically 30-50M parameters for ResNet or ViT backbones plus temporal modeling modules).

**Training efficiency.** Despite using a 7B VLLM, LinkedOut’s training is more efficient than end-to-end video encoder training because: (1) frozen VLLM parameters require no gradient computation or optimizer states; (2) feature extraction is performed once offline, not repeated every epoch; (3) only lightweight MoE fusion and recommendation head are trained online.

## 9.4. Amortized Cost Analysis

A key advantage of LinkedOut’s Store-and-Retrieve architecture is that VLLM feature extraction is a one-time offline cost, amortized across all training epochs, hyperparameter searches, and future model updates. In contrast, end-to-end VideoRec methods must re-encode all videos in every training epoch. Over 100 training epochs, this means end-to-end methods perform video encoding 100 $\times$  more frequently than LinkedOut, making our approach 2-5 $\times$  more cost-efficient in total computational budget despite the upfront VLLM cost. Beyond training efficiency, the one-time feature extraction enables: (1) rapid experimentation with different MoE architectures or recommendation modules without re-encoding videos; (2) easy ensemble of multiple recommendation models using shared features; (3) incremental updates for new videos without retraining the entire system; (4) feature reusability across multiple downstream tasks (recommendation, search, content moderation).

## 9.5. Computational Complexity Analysis

We analyze the asymptotic time complexity of LinkedOut’s key operations with respect to:  $N$  (videos),  $U$  (users),  $T$  (sequence length),  $L$  (layers),  $K$  (tokens per layer),  $d$  (feature dimension). Offline feature extraction per video is  $O(d^2)$  (VLLM forward pass), totaling  $O(N \cdot d^2)$  for the entire catalog (one-time cost). Online serving per query is  $O(L \cdot K \cdot d)$  for MoE gating and fusion, plus  $O(N \log N)$  for top-K candidate ranking. Training per epoch is  $O(U \cdot T \cdot L \cdot K \cdot d)$  for MoE module updates. Critically, unlike end-to-end VideoRec methods with per-epoch complexity  $O(U \cdot T \cdot N \cdot d^2)$ ,

LinkedOut’s training complexity is  $O(U \cdot T \cdot L \cdot K \cdot d)$  where  $L \cdot K \ll N$  (6 layers  $\times$  64 tokens vs. thousands of videos), resulting in significantly faster training convergence.

## 9.6. Architectural Advantages

Beyond raw computational metrics, LinkedOut’s Store-and-Retrieve architecture offers several system-level advantages. First, decoupled scaling: feature extraction (GPU-intensive) and recommendation serving (lightweight) can be scaled independently, optimizing resource allocation. Second, model flexibility: the recommendation module can be retrained or A/B tested without re-extracting video features, enabling rapid experimentation and online learning. Third, cold-start handling: new videos can be recommended immediately after feature extraction, without waiting for collaborative filtering signals to accumulate. Fourth, multi-task support: extracted features support multiple downstream applications (recommendation, search, duplicate detection, content moderation) without redundant encoding. These architectural benefits make LinkedOut practical for production deployment, where system flexibility and maintainability are as important as raw performance metrics.

## 10. Discussion

### 10.1. Mid-Level Layer Dominance

Our layer-wise contribution analysis reveals a surprising finding: intermediate layer L8 contributes 40.9% on average, significantly outperforming both early layers (L0: 16.4%) and late layers (L20: 20.5%). This challenges the common assumption in vision-language models that the deepest layers contain the most task-relevant knowledge. We hypothesize three potential explanations for this phenomenon. First, video recommendation requires balancing fine-grained visual details (e.g., specific objects, scenes, aesthetics) with abstract semantic concepts (e.g., genre, mood, topic). Early layers (L0, L4) primarily capture low-level visual features such as edges, textures, and local patterns, which may be too specific for high-level preference modeling. Late layers (L16, L20) encode abstract reasoning and global scene understanding, which may be over-abstracted for distinguishing similar videos. Mid-level layers (L8, L12) strike an optimal balance by capturing cross-modal alignments, object parts, and semantic groupings that are both specific enough to distinguish content and general enough to generalize across user preferences. Second, video LLMs are pre-trained to align visual and textual concepts across multiple semantic levels. Mid-level layers are where this cross-modal alignment is most active, as shown by prior interpretability studies [15, 41]. For recommendation, the ability to map visual content to semantic tags (e.g., “#vegan\_cooking”, “#urban\_photography”) requires precisely this type of cross-

modal grounding, which is most pronounced in intermediate layers. Third, while video LLMs are trained for language generation, video recommendation prioritizes discriminative visual understanding rather than generative reasoning. The final layers of LLMs are optimized for next-token prediction and causal reasoning, which may not be directly transferable to ranking-based recommendation. In contrast, mid-level representations retain richer perceptual information that can be more effectively adapted to preference matching through our lightweight MoE fusion.

### 10.2. Content-Adaptive Layer Activation Patterns

Beyond aggregate statistics, our sample-level heatmap analysis (Figure 1) reveals significant variation in layer activation patterns across different videos. This adaptive behavior suggests that LinkedOut dynamically selects semantic levels based on video content characteristics. Videos with rich visual details (e.g., cooking tutorials, product reviews) tend to activate L0 more strongly, suggesting that fine-grained visual features are critical for these content types. Conversely, videos with abstract or conceptual content (e.g., educational talks, vlogs) emphasize L20, indicating that high-level reasoning is more valuable. This content-aware adaptation demonstrates that our MoE gating network learns meaningful patterns beyond simple averaging. The ability to leverage different semantic levels for different video types has important implications for cold-start recommendation. New videos with limited interaction data can still be effectively represented by extracting appropriate semantic features from the video LLM. For example, visually rich videos can rely on early-layer features even without textual metadata, while abstract videos can leverage late-layer conceptual understanding. The observed layer specialization suggests that LinkedOut may generalize well to other video domains beyond micro-videos. For instance, movie recommendation might benefit more from late-layer features (capturing narrative and thematic elements), while fashion or product videos might prioritize early-layer features (capturing visual aesthetics and details). Future work could explore domain-specific layer selection strategies.