

Supplementary Material

001 1. Additional Experimental Details

002 1.1. Baseline Architectural Summary

003 For completeness, we provide implementation-level details
004 of the CAV-MAE Sync backbone that are not fully de-
005 scribed in the main paper. These clarify how the baseline
006 organizes tokens, routes information, and separates objec-
007 tives during training.

008 **Token Organization.** For each modality, CAV-MAE
009 Sync prepends: (i) a *dedicated global token* used exclu-
010 sively for the contrastive pathway, and (ii) a set of *regis-*
011 *ter tokens* that act as intermediate buffers within the trans-
012 former. These tokens do not participate in patch recon-
013 struction. Instead, they stabilize the attention flow and pre-
014 vent patch tokens from absorbing semantic responsibilities
015 needed for contrastive learning.

016 **Joint-Layer Execution Structure.** Unlike the original
017 CAV-MAE, Sync performs *three forward passes* through
018 the joint transformer layer J , each with its own LayerNorm:
019 • visual-only pass: $J_v(Z^v)$,
020 • audio-only pass: $J_a(Z^a)$,
021 • fused pass for reconstruction: $J_{\text{fuse}}([\text{vis}(Z^v); \text{vis}(Z^a)])$.
022 The first two passes generate modality-specific global rep-
023 resentations *without* involving masked tokens. The fused
024 pass is used only by the reconstruction decoder.

025 **Global Token Pathway.** The global tokens are not aver-
026 aged from patches (as in MAE and CAV-MAE). Instead,
027 after the single-modality encoders, the global token is fed
028 directly into its corresponding joint-layer pass:

$$029 \quad g^v = \text{LN}_v(J_v(Z^v)_0), \quad g^a = \text{LN}_a(J_a(Z^a)_0),$$

030 where the subscript 0 denotes the global-token position.
031 This ensures that contrastive learning operates on explic-
032 itly learned semantic tokens rather than statistics of visible
033 patches.

034 **Reconstruction Pathway.** Only *visible patch tokens* are
035 forwarded to J_{fuse} . Global tokens and register tokens are
036 *removed* from the reconstruction stream. This prevents gra-
037 dients from the generative objective from flowing through
038 global tokens and interfering with cross-modal alignment.
039 The joint decoder D receives a concatenation of recon-
040 structed visual and audio features and predicts masked
041 patches for *both* modalities.

Masking Semantics. CAV-MAE Sync applies independ- 042
ent random masks to visual and audio patches. Because 043
global tokens do not depend on masked inputs, their gradi- 044
ents are unaffected by the stochastic patch visibility, reduc- 045
ing semantic noise compared to CAV-MAE where global 046
descriptors are aggregated from partially visible patches. 047

Temporal Handling. A key difference from the standard 048
CAV-MAE baseline is that Sync treats audio as a sequence 049
aligned to the T sampled video frames. The baseline there- 050
fore generates: 051

- T audio segments (one per frame), 052
- T audio global tokens, 053
- T visual global tokens, 054

enabling multi-token temporal contrast without modifying 055
the backbone architecture. 056

057 1.2. Data Preprocessing

We follow the official CAV-MAE Sync preprocessing 058
pipeline without modification. Each video is uniformly 059
sampled into $T=16$ frames, which are resized and center- 060
cropped to 224×224 before being patchified into 16×16 vi- 061
sual patches. The corresponding audio waveform is con- 062
verted into 128-bin log-Mel filterbanks using a 25 ms Han- 063
ning window and 10 ms hop. For every sampled frame, 064
a temporally aligned 4-second spectrogram segment (size 065
 128×416) is extracted by mapping the frame index to its 066
spectrogram coordinate, ensuring strict frame-level tempo- 067
ral correspondence. These audio segments are then patchi- 068
fied with the same 16×16 grid, producing 208 audio to- 069
kens per segment. Standard normalization is applied to both 070
modalities, and the resulting sequences serve as the input 071
token streams for subsequent masking and encoding in the 072
baseline model. 073

In practice, to reduce I/O overhead during large-scale 074
pretraining, we additionally package all audio–visual sam- 075
ples into `webdataset` shards. Each shard stores paired 076
RGB frames, their aligned audio segments, and the cor- 077
responding metadata, enabling sequential and locality- 078
friendly reads during distributed training. This packag- 079
ing significantly alleviates random-access bottlenecks in- 080
herent to filesystem-based loading, allowing the dataloader 081
to stream pre-compressed samples efficiently across multi- 082
ple GPUs and improving overall training throughput. 083

Setting	Pretrain	AS20K LP	VGG LP
Dataset	AS-2M	AS-20K	VGGSound
Optimizer	Adam	Adam	Adam
Learning Rate	2×10^{-4}	5×10^{-2}	1×10^{-3}
LR Scheduler	Cosine	Cosine	Cosine
Epochs	35	15	10
Warmup Epochs	3.5	1.5	1
Batch Size	8×64	48	48
GPUs	8×H100	2×H100	2×H100
Audio Input Size	128×416	16×128×416	16×128×416
Class-Balanced	No	No	Yes
Mixup	No	Yes	Yes
Random Shift	Yes	Yes	Yes
Loss Function	—	BCE	CE
Weight Averaging	No	Yes	Yes
Norm Mean	-5.081	-5.081	-5.081
Norm STD	4.485	4.485	4.485

Table 1. Hyperparameters for pretraining and linear probing.

1.3. Pretraining Configuration

Our pretraining setup follows the official CAV-MAE Sync configuration, including optimizer settings, masking ratio, data augmentations, batch size, and loss weights. All experiments are conducted on $8 \times$ NVIDIA H100 GPUs. Due to missing entries in our local AudioSet-2M copy, we extend the pretraining schedule from 25 to 35 epochs, and proportionally adjust warmup and cosine learning-rate scheduling. Refer to 1 for detailed configurations.

Although our dual-path framework introduces an additional teacher forward pass, the extra computation is modest because the teacher processes only short sequences. In practice, the per-epoch time increases from 730 s to 1045 s, resulting in a total pretraining time of approximately 7.1 h. Table 2 summarizes this overhead.

Model	Per-epoch (s)	Total
baseline	730	7.1 h
Ours	1045	10.2 h

Table 2. Pretraining runtime comparison. The dual-path teacher branch increases cost moderately while remaining efficient.

1.4. Downstream Evaluation

For all downstream tasks, we keep the implementation consistent with the original CAV-MAE Sync settings to ensure comparability. This includes the same optimizer (Adam with $(\beta_1, \beta_2) = (0.95, 0.999)$), cosine learning-rate schedule with warmup, masking strategy, and lightweight data augmentations applied during linear probing. Unless otherwise specified, all encoder weights are frozen and only the task-specific heads are trained.

Zero-shot Retrieval. Since retrieval is not fully detailed in the main paper, we summarize the procedure here. For each video, the model produces a temporal sequence of visual and audio global tokens, $\{g_t^v\}_{t=1}^T$ and $\{g_t^a\}_{t=1}^T$. For any query-candidate pair, a $T \times T$ cosine-similarity matrix is computed between these sequences. The final similarity score is obtained by averaging the diagonal elements, which corresponds to comparing frame-aligned audio-visual tokens. This score is used to rank candidates, and $\text{Recall}@\{1, 5, 10\}$ is reported for both $V \rightarrow A$ and $A \rightarrow V$ retrieval on the standard subsampled splits.

Classification. For AudioSet-20K and VGGSound classification, we train only the two-layer classification head while keeping the pretrained encoders fixed. We use binary cross-entropy for AudioSet and standard cross-entropy for VGGSound. The same sampling strategy, data normalization, and mixup settings as in the original configuration are used. The learning rate follows a cosine schedule with task-dependent peak values.

2. Visualization and Qualitative Analysis

2.1. Training Dynamics Analysis

Figure 1 compares the training behavior of the original CAV-MAE Sync baseline with our dual-path formulation. We report the per-step reconstruction loss and contrastive loss throughout pretraining. A clear pattern emerges: in both models, the reconstruction objective converges extremely quickly, stabilizing within the first few hundred iterations. In contrast, the contrastive objective decreases at a significantly slower rate and continues to improve for thousands of steps. This divergence in optimization speed suggests that the generative and discriminative objectives are only weakly coupled during training and naturally evolve on different timescales.

This observation supports the rationale behind our decoupling strategy: since reconstruction and alignment exhibit largely independent behaviors, optimizing them in separate pathways is a principled choice. Moreover, the dual-path design accelerates the optimization of the contrastive branch, as shown by the consistently faster decay of the contrastive loss in our model. By preventing gradient interference and ensuring both pathways receive clean, unconflicted supervision signals, the dual-path formulation enables more stable optimization and faster convergence, particularly for the alignment objective that directly drives cross-modal semantic learning.

We additionally observe a consistent gap between the reconstruction losses of the two modalities: both the baseline and our dual-path model stabilize around a vision reconstruction loss of approximately 0.35, whereas the audio reconstruction loss remains notably higher at roughly 0.60.

160 This discrepancy indicates that the model finds it substan-
161 tially easier to encode and reconstruct visual tokens than
162 audio tokens, reflecting the stronger spatial structure and
163 inductive biases inherent to vision Transformers. As dis-
164 cussed in the main paper, this imbalance suggests that cross-
165 modal interaction is often visually anchored, with the model
166 more readily leveraging visual representations during align-
167 ment and downstream tasks. The consistent reconstruction
168 gap highlights an intrinsic asymmetry between modalities
169 and further motivates architectural designs that prevent the
170 visually dominant representations from overwhelming the
171 audio pathway during joint learning.

172 2.2. Cross-Modal Retrieval Case Study

173 To better understand the qualitative advantages of our ap-
174 proach, Figure 3 presents a cross-modal retrieval case study
175 comparing the top retrieved videos from the baseline CAV-
176 MAE Sync and our method under various audio queries.
177 Each row corresponds to a distinct sound event, including
178 *roller coaster running*, *people shuffling*, *basketball bounce*,
179 and *dog barking*.

180 Across all examples, the baseline often returns videos
181 that share only coarse or incidental correlations with the
182 audio query (e.g., unrelated scenes with background noise
183 or visually mismatched contexts). This behavior reflects
184 the limitations of global audio representations and coupled
185 reconstruction–contrastive optimization, which tend to di-
186 lute fine-grained semantic cues necessary for discrimina-
187 tive retrieval. In contrast, our method consistently retrieves
188 videos whose visual content closely matches the semantics
189 and temporal characteristics of the input audio.

190 2.3. Global Token Embedding Visualization

191 To qualitatively assess the semantic structure of the learned
192 representations, we project the audio and video global to-
193 kens into a shared 2D space using t-SNE. We randomly
194 sample 10 classes from the VGGSound test split and en-
195 code both modalities with the pretrained models. Audio
196 and video embeddings are visualized using distinct marker
197 shapes (circles for audio, triangles for video), while colors
198 indicate different semantic categories.

199 As shown in Fig. 2, the baseline model exhibits notice-
200 able modality discrepancy: audio and video embeddings
201 belonging to the same class are distributed across distant
202 regions, and the resulting clusters often display weak se-
203 mantic coherence. This misalignment reflects the inherent
204 limitations of joint reconstruction–alignment optimization,
205 where random masking and inconsistent gradients introduce
206 semantic noise.

207 In contrast, our method produces markedly tighter struc-
208 tures. Audio and video embeddings of the same class
209 form compact, well-overlapped co-clusters, while clus-
210 ters of different classes become more clearly separated.

211 These improvements confirm that our teacher-guided dual-
212 path framework successfully enhances global semantic
213 alignment, enabling the model to learn a more modality-
214 consistent and discriminative embedding space.

211
212
213
214

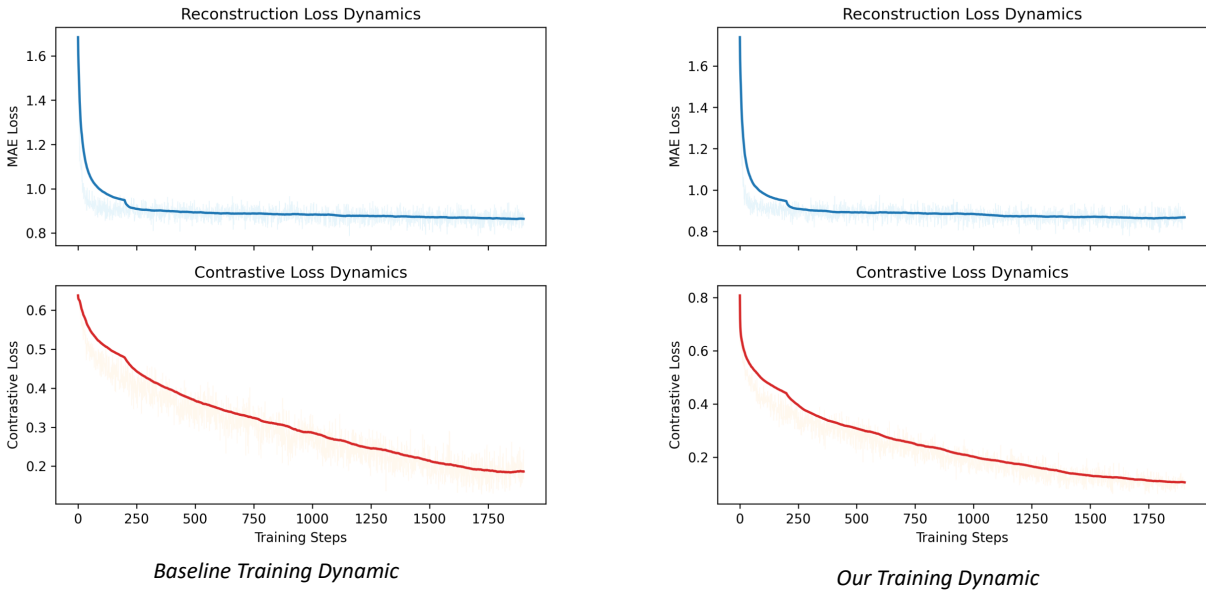


Figure 1. Training loss dynamics of the baseline CAV-MAE Sync (left) and our dual-path formulation (right). For both models, the reconstruction loss quickly converges within the first few hundred steps, while the contrastive loss decreases much more slowly, indicating that the two objectives behave largely independently during optimization. Our dual-path strategy further accelerates the convergence of the contrastive branch, supporting the motivation of separating reconstruction and alignment pathways for more stable and efficient training.

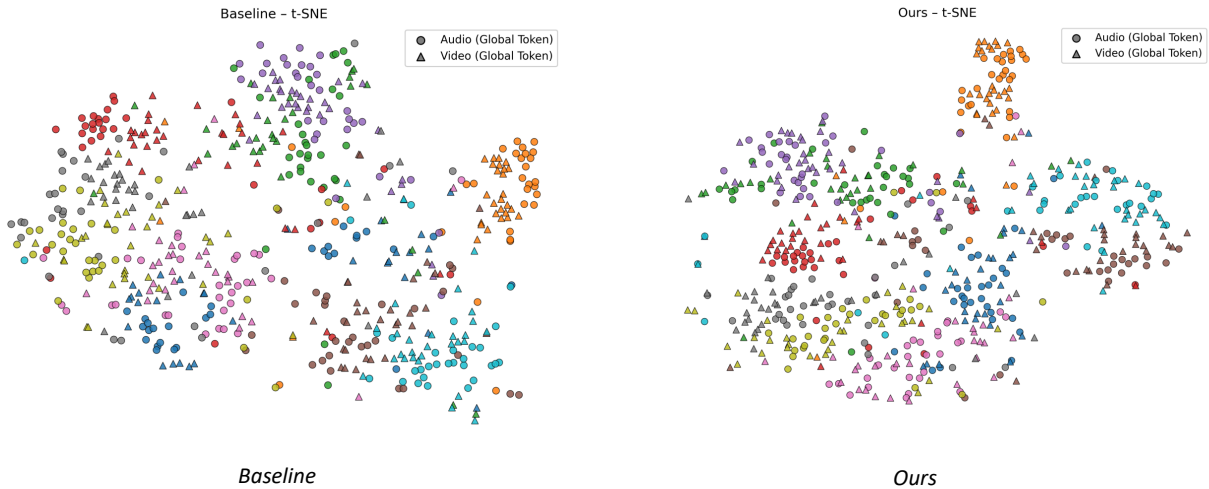


Figure 2. t-SNE visualization of global token embeddings from the baseline (left) and our method (right). Each point denotes a global embedding, with **circles for audio** and **triangles for video**, and colors indicating different semantic classes. The baseline exhibits modality discrepancy compared to ours. Audio and video clusters of the some class are loosely aligned and often drift apart. In contrast, our method produces **tight audio–video co-clusters** with clearer inter-class separation, demonstrating substantially improved cross-modal consistency learned during pretraining.

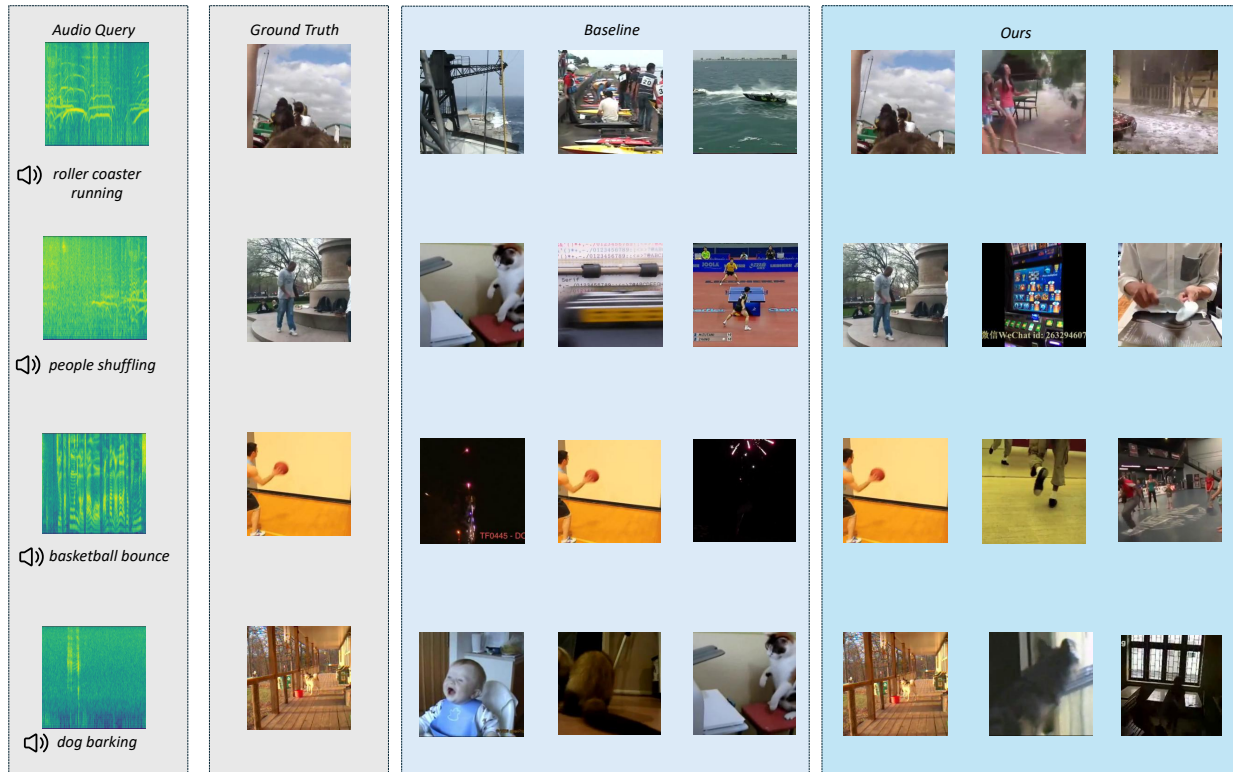


Figure 3. **Qualitative cross-modal retrieval results.** Given an audio query (left), we compare the top retrieved videos from the **baseline CAV-MAE Sync** (middle) and **our method** (right). Each row corresponds to a different audio event (*roller coaster running*, *people shuffling*, *basketball bounce*, *dog barking*). While the baseline frequently retrieves semantically irrelevant or visually mismatched scenes, our model consistently returns videos that better match both the *sound semantics* and the *visual content*, demonstrating substantially improved audio–visual alignment learned during pretraining.