Ap	opendix
A	SUMMARY
The	appendix is organized as follows:
	• § B detailed training and evaluation settings of our models including hyper-parame regarding models and optimizers.
	• § C presents a comprehensive introduction on the datasets we use for evaluation and t corresponding metrics.
В	TRAINING CONFIGURATION
<b>B</b> .1	PRETRAINING SETTINGS
We corp imp data has furtl	detail the specific pretraining configurations of MiCo, focusing on the multi-dataset joint train pora, the dataset mix ratios for each corpus, and the learning objectives for each corpus. rove data quality, we employed a trained vision captioner to generate new captions for the CC sets, replacing the original captions. Although MiCo has only been trained for 300,000 step already demonstrated outstanding performance on various downstream tasks. We anticipate her increasing the number of training steps will significantly enhance the model's capabilitie
The the the	pretraining of MiCo involves a combination of different datasets, each contributing uniquel model's learning process. With a parameter size of 1.0 billion and a sample size of 334 mill model utilizes a diverse training corpus to achieve its results.
1. V of 2	AST-27M: This dataset contributes 324 million samples to the training process. With a batch 048, the model undergoes 160,000 steps, completing one epoch.
2. V spai	ALOR-1M: In this dataset, 1 million samples are used with a batch size of 1024. The train as 70,000 steps, which equates to approximately 71.7 epochs.
3. V in to resu	VavCaps, CC4M, and WebVid-2.5M: These datasets are combined, contributing 9 million sam otal. The batch size for this combined dataset is 1024, and the model is trained over 70,000 st liting in 8.0 epochs.
The qual repr	careful selection and combination of these datasets, along with the application of new, h lity captions for the CC4M datasets, enhance the training efficiency and the quality of the lear resentations.
B.2	FINE-TUNING SETTINGS
We traiı fran	detail the downstream task finetuning settings, specifying the learning rate, batch size, ep ning objectives, and resolution. The configurations also include the number of sampled vin nes or audio clips used in training and testing phases. Here are the comprehensive settings:
Ret	rieval Tasks (RET)
	Image-Text Modality
	- MSCOCO: Learning rate of 1e-5, batch size of 256, 5 epochs, with the objective
	retrieval, and a resolution of 384.
	- Flickr: Learning rate of 1e-5, batch size of 256, 5 epochs, with the objective retrieval and a resolution of 384
	Audio_Text Modelity (A-T)
	Clothe V1/V2. Learning rate of 20.5. herebeing of 64.10 encodes with the object
	- CIOLIDO V 1/ V 2: Learning rate of 2e-5, batch size of 64, 10 epochs, with the object

for retrieval, using 3 audio clips during both training and testing.

918	- AudioCaps: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective for
919	retrieval, using 1 audio clip during both training and testing.
920	• Multi-modal (MM)
921	- MSPVTT: Learning rate of 2e 5 batch size of 64.3.6 enochs, with the objective for
922	retrieval using 8 video frames during training and 16 during testing, with a resolution
923	of 224.
924	- YouCook2: Learning rate of 3e-5, batch size of 64, 30 epochs, with the objective for
920	retrieval, using 8 video frames during training and 16 during testing, with a resolution
920	of 224.
928	- VALOR-32K: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective for
929	retrieval, using 8 video frames during both training and testing, with a resolution of
930	224.
931	- VATEX: Learning rate of 2e-5, batch size of 64, 2.5 epochs, with the objective for
932	retrieval, using 8 video frames during training and 16 during testing, with a resolution
933	of 224.
934	- <b>DiDeMo</b> : Learning rate of 2e-5, batch size of 64, 40 epochs, with the objective for
935	during both training and testing, with a resolution of 224
936	ANET: Learning rate of 20.5, botch size of 64, 20 another with the objective for
937	- AINE 1. Learning rate of 2e-3, batch size of 04, 20 epochs, with the objective for retrieval using 8 video frames during training and 32 during testing, and 2 audio clins
938	during both training and testing, with a resolution of 224.
939	
940	Captioning Tasks (CAP)
941	
942	• Image-Text Modality
943	- MSCOCO: Learning rate of 1e-5, batch size of 64, 5 epochs, with the objective for
944	caption, and a resolution of 480.
945	- MSCOCO(SCST): Learning rate of 2.5e-6, batch size of 64, 2.5 epochs, with the
940	objective for caption, and a resolution of 480.
948	• Audio-Text Modality (A-T)
949	- ClothoV1/V2: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective
950	for caption, using 3 audio clips during both training and testing.
951	- AudioCaps: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective for
952	caption, using 1 audio clip during both training and testing.
953	• Multi-modal (MM)
954	- MSRVTT: Learning rate of 2e-5, batch size of 128, 10 epochs, with the objective for
955	caption, using 8 video frames during both training and testing, with a resolution of 224.
956	- YouCook2: Learning rate of 3e-5, batch size of 64, 30 epochs, with the objective for
957	caption, using 8 video frames during training and 16 during testing, with a resolution
958	of 224.
959	- VALOR-32K: Learning rate of 1e-5, batch size of 64, 10 epochs, with the objective for
960	caption, using 8 video frames during training and 12 during testing, with a resolution
301 060	01 224.
902	Ouestion Answering Tasks (OA)
964	
965	Visual-Text Modality (Vis)
966	- MSVD-QA: Learning rate of 1e-5, batch size of 64, 10 epochs, with the objective for
967	QA, using 8 video frames during training and 14 during testing, with a resolution of
968	224.
969	- TGIF-FrameQA: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective
970	for QA, using 4 video frames during both training and testing, with a resolution of 224.
971	<ul> <li>VQAv2: Learning rate of 2e-5, batch size of 128, 20 epochs, with the objective for QA, and a resolution of 384.</li> </ul>

Table 8: Detailed training configurations of MiCo for multimodal learning. Apart from the configurations
 shown in the table, for image tasks, we use random left-right flipping, random resized crop, color jitter of 0.4,
 Auto-augment, and no repeated augmentation for every model.

975									
070	eettinge	Image		Audio		Video		Depth & Normal Map	
976	settings	ViT-L	ViT-g	ViT-L	ViT-g	ViT-L	ViT-g	ViT-L	ViT-g
977	Input Shape	224	224	224	224	224	224	224	224
978	batch size	4096	512	4096	512	4096	512	4096	512
070	optimizer	AdamW							
979	LR	$4 \times 10^{-3}$	$5 \times 10^{-5}$						
980	LR schedule	cosine							
981	weight decay	0.05	$1 \times 10^{-8}$						
982	warmup epochs	5	0	5	0	5	0	5	0
002	epochs	90	30	90	30	90	20	90	20
903	mixup alpha	0.8	0.0	0.8	0.0	0.8	0.0	0.8	0.0
984	cutmix alpha	1.0	0.0	1.0	0.0	1.0	0.0	1.0	0.0
985	erasing prob.	0.25	0.25	0.25	0.25	0.25	0.25	0.25	0.25
986	dropout rate	0.1	0.2	0.1	0.2	0.1	0.3	0.2	0.3

Algorithm 1 Multimodal Context Pretraining Algorithm, PyTorch-like

```
def train(video_pixels=None, image_pixels=None, depth_pixels=None, audio_spectrograms=
   None):
  # Get Mixed Data
  modal inputs = [video pixels, image pixels, depth pixels, audio spectrograms]
  modal captions = [video captions, image captions, depth captions, audio captions]
  # Extract Features
  modal feats = [self.encoder(modal) for modal in modal inputs if modal is not None]
  multimodal feats = torch.cat(modal feats)
  concatenated captions = ''.join(modal captions)
  text_feats = self.text_encoder(concatenated_captions)
  # Losses
  contra loss = Contrasive Loss(multimodal feats, text feats)
  matching_loss = Matching_Loss(modal_captions, multimodal feats)
  gen loss = Generation Loss(modal captions.mask(0.6), multimodal feats)
  # Total Loss
  loss = contra_loss + matching_loss + gen_loss
return loss
```

• Multi-modal (MM)

987

988 989

990

991

992

993 994

995

996

997

998 999

1000

1001

1002

1004

1008 1009

1010 1011

1012

1013

1014

1015

1016

1017

1020

- MSRVTT-QA: Learning rate of 2e-5, batch size of 64, 4.5 epochs, with the objective for QA, using 8 video frames and 1 audio clip during both training and testing, with a resolution of 224.
- MUSIC-AVQA: Learning rate of 2e-5, batch size of 64, 20 epochs, with the objective for QA, using 8 video frames and 2 audio clips during both training and testing, with a resolution of 224.
- ANET-QA: Learning rate of 2e-5, batch size of 64, 10 epochs, with the objective for QA, using 8 video frames during training and 16 during testing, and 2 audio clips during both training and testing, with a resolution of 224.
- 1021 These settings have been optimized to balance efficiency and performance, even though most hyperparameters are not precisely tuned.
- For evaluation purposes, we employ different strategies tailored to specific tasks:
- 1025 1. Retrieval Tasks: All candidates are initially ranked using Omni-modal Contrastive Loss. Following this, the Top-50 candidates undergo a reranking process through the Omni-modal Matching Process.

```
1026
       Algorithm 2 Dataset Split Algorithm
1027
       import pandas as pd
1028
       from sklearn.model selection import train test split
1029
1030
       # Assume 'data' is a DataFrame containing the full dataset with columns ['category',
1031
            vision caption', 'audio caption', 'depth', 'normal', 'subtitle']
       # Adding an 'index' column to keep track of the original indices
1032
       data['index'] = data.index
1033
1034
       # Define the sizes of each subset
1035
       subset sizes = [1e6, 1e7, 1.1e7, 3.34e7]
1036
1037
        # Function to create stratified samples
       def create_subset(data, size):
          subset, _ = train_test_split(data, train_size=size, stratify=data['category'],
1039
               random state=42)
1040
           return subset
1041
       # Creating subsets
       subset 1M = create subset(data, 1e6)
1043
       subset 10M = create subset(data, 1e7)
1044
       subset 110M = create subset(data, 1.1e7)
1045
       subset 334M = create subset(data, 3.34e7)
1046
1047
       # Reset index for each subset
       subset 1M.reset index(drop=True, inplace=True)
1048
       subset 10M.reset index(drop=True, inplace=True)
1049
       subset_110M.reset_index(drop=True, inplace=True)
1050
       subset 334M.reset index(drop=True, inplace=True)
1051
```

1052 1053

2. Captioning Tasks: Beam search with a beam size of 3 is utilized to generate captions, ensuring a comprehensive exploration of possible outputs.

1056 3. Question Answering (QA) Tasks: These are treated as open-ended generative problems. Questions are used as prefixes, and answers are generated without any constraints, allowing for flexible and contextually appropriate responses.

For comparisons with state-of-the-art (SOTA) models and ablation studies, we use the following evaluation metrics: 1) Retrieval Tasks: Recall@1. 2) Captioning Tasks: CIDEr. 3) QA Tasks:
Accuracy (Acc) These metrics provide a comprehensive assessment of the model's performance across different types of tasks.

1063 1064

C DATASETS AND METRICS

- 1066
- 1067

1068 Dataset Split To split the mix of datasets into subsets of 1M, 10M, 110M, and 334M video clips while preserving its diversity and quality, we employed a proportional stratified sampling method. Initially, 1069 the dataset, which spans over 15 categories (including music, gaming, education, entertainment, and 1070 animals) and includes vision, audio, depth, normal maps, and text modalities, was organized and 1071 labeled. Stratified random sampling was then used to ensure each subset accurately reflected the 1072 distribution of categories and modalities present in the full dataset. This method involved selecting samples proportionally from each category to maintain representative distributions. The vision and 1074 audio captions were also kept proportional in length and quantity, ensuring that each subset retained 1075 the comprehensive characteristics of the original dataset.

1076 1077

1078

C.1 SINGLE-MODALITY EVALUATION DETAILS

**Text**. The MMLU (Massive Multitask Language Understanding) benchmark is designed to evaluate the multitask accuracy of language models across 57 diverse tasks, including subjects like mathemat-

ics, history, and biology. It assesses models' abilities to generalize and apply knowledge in various domains, providing a comprehensive measure of text understanding and reasoning skills.

Image. We conduct experiments on ImageNet-1K (Deng et al., 2009), a dataset comprising approximately 1.3 million images across 1,000 categories. In line with common practices (Wang et al., 2021a; Liu et al., 2021; 2022; Ding et al., 2023), base-scale models are trained for 300 epochs. Large-scale models undergo pre-training on ImageNet-22K, which includes 14.2 million images, for 90 epochs, followed by fine-tuning on ImageNet-1K for an additional 20 epochs.

**Thermal and Hyperspectral data understanding**. We conduct experiments on infrared image recognition using the RegDB dataset, X-ray scan analysis with the Chest X-Ray dataset (Rahman et al., 2020), and hyperspectral data recognition using the Indian Pine dataset<sup>2</sup>.

1091 Depth. The NYU Depth Dataset (NYU-D) comprises RGB and depth image pairs captured from indoor scenes. It includes 1,449 densely labeled pairs for training and testing, along with over 400,000 unlabeled frames.

Audio. For audio recognition, Audioset-2M dataset comprises over 2 million human-labeled 10 second audio clips drawn from YouTube videos. It covers a wide range of 527 sound event classes,
 providing a comprehensive resource for training and evaluating audio event detection and classifica tion models.

Video. The Kinetics-700 dataset contains 700,000 video clips covering 700 human action classes, used for action recognition tasks. The MSR-VTT dataset includes 10,000 video clips paired with multiple textual descriptions, supporting video captioning, retrieval, and content understanding research.

Time-series. Global Weather Forecasting (Wu et al., 2023) includes global, regional, and Olympics data from NCEI and CMA, comprising hourly weather measurements from thousands of stations. Evaluation involved splitting data into training, validation, and test sets (7:1:2) using MSE and MAE metrics.

Graph. PCQM4M-LSC dataset is a large-scale collection of 4.4 million organic molecules, each with up to 23 heavy atoms and associated quantum-mechanical properties. Aimed at predicting molecular properties through machine learning, this dataset is highly relevant for applications in drug discovery and material science.

**Tabular**. The fraud dataset comprises transaction records, including features like transaction amount,
location, time, and user information. It is designed for machine learning models to detect fraudulent
activities. This dataset is crucial for developing and testing algorithms to enhance security in financial
systems and reduce economic losses due to fraud.

1115 IMU. The Ego4D dataset includes inertial measurement unit (IMU) data captured from wearable
 1116 devices, providing detailed motion and orientation information. This dataset supports research in
 1117 human activity recognition, augmented reality, and robotics, offering comprehensive insights into
 1118 human movements and interactions with the environment.

- 1119
- 1120
- 1121

## 1122 C.2 CROSS-MODALITY EVALUATION DETAILS

1123 1124

We evaluated MiCo across several well-known downstream datasets, including MSRVTT, VATEX,
 YouCook2, VALOR-32K, MSVD, DiDeMo, ActivityNet Caption, TGIF, MUSIC-AVQA, Clotho,
 AudioCaps, MSCOCO, Flickr30K, and VQAv2. The specific train/validation/test splits for these
 benchmarks are detailed below:

1131

<sup>1129</sup> 

<sup>1130</sup> 

<sup>1132</sup> 

<sup>1133 &</sup>lt;sup>2</sup>https://github.com/danfenghong/IEEE\_TGRS\_SpectralFormer/blob/main/data/IndianPine.
mat

1134 1135	Retrieval Tasks
1136	Audio-Text Modality (A-T)
1138 1139 1140	• <b>ClothoV1</b> (Drossos et al., 2020): This dataset includes 2,893 audio clips for training and 1,045 for validation. The corresponding captions number 14,465 for training and 5,225 for validation.
1141 1142	• <b>ClothoV2</b> (Drossos et al., 2020): Contains 3,839 audio clips for training and 1,045 for validation, with 19,195 captions for training and 5,225 for validation.
1143 1144 1145 1146	• AudioCaps (Kim et al., 2019): Comprises 49,291 audio clips for training, 428 for validation, and 816 for testing, along with 49,291 captions for training, 2,140 for validation, and 4,080 for testing.
1147 1148	VIDEO-TEXT MODALITY (V-T)
1149 1150 1151 1152	<ul> <li>MSRVTT (Xu et al., 2016a): Comprises 10K video clips and 200K captions, spanning diverse topics such as human activities, sports, and natural landscapes. We evaluate text-to-video retrieval, video captioning, and video QA using this dataset. Contains 9,000 videos for training and 1,000 for testing, with 180,000 captions for training and 1,000 for testing.</li> </ul>
1153 1154 1155 1156 1157	• YouCook2 (Zhou et al., 2018): Comprises 14K video clips extracted from 2K instructional cooking videos on YouTube. Each video features multiple actions performed by chefs, along with corresponding textual descriptions and temporal annotations. Includes 10,337 videos for training and 3,492 for validation, with matching captions.
1158 1159 1160 1161 1162	• VALOR-32K (Chen et al., 2023a): An audiovisual video-language benchmark containing 32K 10-second video clips sourced from AudioSet (Gemmeke et al., 2017). Each clip includes annotations with captions that describe both the visual and audio content. Consists of 25,000 videos for training, 3,500 for validation, and 3,500 for testing, each with corresponding captions.
1163 1164 1165 1166 1167	• <b>DiDeMo</b> (Anne Hendricks et al., 2017): Comprises 10K long-form videos sourced from Flickr, with each video annotated with four short sentences in temporal order. For this benchmark, we concatenate these short sentences and evaluate 'paragraph-to-video' retrieval, using the official split. Features 8,394 videos for training, 1,065 for validation, and 1,003 for testing, along with their captions.
1168 1169 1170 1171	• ActivityNet (ANET) (Krishna et al., 2017): Includes 20K long-form videos (average length of 180 seconds) from YouTube, accompanied by 100K captions. We evaluate text-to-video retrieval and video QA on this dataset. Comprises 10,009 videos for training and 4,917 for testing, with corresponding captions.
1172 1173 1174 1175	• LSMDC (Rohrbach et al., 2017): Contains 101,046 videos for training, 7,408 for validation, and 1,000 for testing, with corresponding captions.
1176 1177	CAPTIONING TASKS
1178 1179	Audio-Text Modality (A-T)
1180 1181 1182	• <b>ClothoV1</b> (Drossos et al., 2020): This dataset includes 2,893 audio clips for training and 1,045 for validation. The corresponding captions number 14,465 for training and 5,225 for validation.
1183 1184	• <b>ClothoV2</b> (Drossos et al., 2020): Contains 3,839 audio clips for training and 1,045 for validation, with 19,195 captions for training and 5,225 for validation.
1186 1187	• AudioCaps (Kim et al., 2019): Comprises 49,838 audio clips for training, 495 for validation, and 975 for testing, along with 49,438 captions for training, 2,475 for validation, and 4,875 for testing.

VIDEO-TEXT MODALITY (V-T)
• MSRVTT (Xu et al., 2016a): Contains 6.513 videos for training, 497 for validation, and
2,990 for testing, with 130,260 captions for training, 9,940 for validation, and 59,800 for
testing.
• YouCook2 (Zhou et al., 2018): Includes 10.337 videos for training and 3.492 for validation.
with matching captions.
• VALOR-32K (Chen et al., 2023a): Consists of 25,000 videos for training, 3,500 for validation, and 3,500 for testing, each with corresponding captions.
• VATEX (Wang et al., 2019): Consists of 41,250 video clips sourced from the Kinetics-600 dataset (Kay et al., 2017), accompanied by 825,000 sentence-level descriptions. Contains 25,991 videos for training, 3,000 for validation, and 6,000 for testing, with 259,910 captions for training, 30,000 for validation, and 60,000 for testing.
QUESTION ANSWERING (QA) TASKS
VIDEO-TEXT MODALITY (V-T)
• MSRVTT-QA (Xu et al., 2017): Contains 6,513 videos for training, 497 for validation, and 2,990 for testing, with 158,581 QA pairs for training, 12,278 for validation, and 72,821 for testing.
• <b>MUSIC-AVQA</b> (Li et al., 2022): An audiovisual video QA benchmark containing over 45K Q-A pairs, covering 33 different question templates across various modalities and question types. Includes 9,277 videos for training, 3,815 for validation, and 6,399 for testing, with 32,087 QA pairs for training, 4,595 for validation, and 9,185 for testing.
• ANET-QA (Yu et al., 2019a): Comprises 3,200 videos for training, 1,800 for validation, and 800 for testing, with 32,000 QA pairs for training, 18,000 for validation, and 8,000 for testing.
IMAGE-BASED TASKS
• <b>MSCOCO</b> (Lin et al., 2014): Comprises 123K images, each paired with 5 annotated captions. We evaluate text-to-image retrieval and image captioning on this dataset.
• Flickr30K (Plummer et al., 2015): Contains 31K images, each paired with five descriptive captions. This dataset is widely used for evaluating image captioning and text-to-image retrieval tasks.
VISUAL QUESTION ANSWERING
• $VOAv2$ (Goval et al. 2017a): A large-scale Visual Question Answering dataset comprising
over 265K images and 1.1M questions, designed to improve the balance of answer types per
question. This dataset is used to evaluate models' abilities to understand and reason about
visual content by providing accurate answers to questions based on the images.