

VISTA: Dense Multi-Label Classroom Coding with Vision-Language Models

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Video-language benchmarks are usually constructed by the*
002 *dataset authors without published reliability statistics, leav-*
003 *ing the noise floor of the construct unknown. We argue that*
004 *multimodal benchmarking benefits from methods taken from*
005 *research communities that have already invested in strate-*
006 *gies to ensure reliability. We illustrate the case with the*
007 *Classroom Observation Protocol for Undergraduate STEM*
008 *(COPUS): a 24-code multi-label observation instrument*
009 *with a decade of peer-reviewed reliability literature. We re-*
010 *cast COPUS as a video benchmark for multimodal founda-*
011 *tion models, where it provides a dense set of structured la-*
012 *bels (a 24-dimensional binary vector every 2 minutes across*
013 *a 50–90 minute lecture), an externally validated vocabu-*
014 *lary, and established literature that provides a per-code re-*
015 *liability target based on human evaluators. Annotations in*
016 *our evaluation corpus are produced by a 5-person human-*
017 *evaluator panel whose consensus matrix is our reference.*
018 *We propose VISTA, a baseline that runs MiniCPM-V-4.5*
019 *over a dense sliding window, refines its per-window out-*
020 *puts with a lightweight multi-layer perceptron (MLP) head*
021 *trained on top of the frozen backbone, and max-pools the*
022 *resulting predictions onto the 2-minute COPUS grid. On*
023 *three held-out chemistry lectures, VISTA reaches 80.1%*
024 *restricted macro accuracy versus 74.9% for the zero-shot*
025 *variant, with the largest residual errors on visually simi-*
026 *lar instructor codes (RTW) and on rare audio-dependent*
027 *codes. We characterise three systematic failure modes*
028 *(audio-partial observability, fine-grained group-work dis-*
029 *crimination, long-tail recall) and release the benchmark*
030 *tooling, prompts and baseline code.*

031 1. Introduction

032 Benchmark renewal in video-language modelling is increas-
033 ingly a data problem. Video-language benchmarks such
034 as Video-MME [3], MVBench [5], LongVideoBench [16],
035 and EgoSchema [7] have spread quickly, but all share struc-
036 tural limitations. First, although several of them use multi-
037 stage annotation pipelines, the ground truth is constructed

entirely by the dataset authors and none report formal inter- 038
rater reliability statistics. The noise floor of the construct 039
is unknown, and there is no externally validated reliability 040
ceiling against which model error can be compared. Sec- 041
ond, the tasks are dominated by short-clip question–answer 042
pairs; dense labelling at regular intervals throughout the 043
video is rare. Third, most benchmarks define their own 044
custom label set, so cross-benchmark scores reflect dataset 045
choices as much as model skill. 046

We argue that the multimodal community underutilises 047
a different kind of benchmark: instruments taken from re- 048
search communities that have spent years validating their 049
reliability. As a concrete case study, we propose the Class- 050
room Observation Protocol for Undergraduate STEM (CO- 051
PUS) [11] as a video benchmark; COPUS is the standard 052
classroom-observation protocol, although other observation 053
protocols such as RTOP [9] and TDOP [4] have also been 054
used. COPUS partitions a 50–90 minute lecture into 2- 055
minute intervals and, for each interval, asks a trained ob- 056
server to mark which of 24 predefined behaviours (13 stu- 057
dent, 11 instructor) occurred. The protocol has been widely 058
used in the STEM education literature, has a published 059
inter-rater reliability record, and is the operational target 060
of education researchers who would benefit from an auto- 061
mated coder. Prior automated analysis of classroom video 062
and audio [2, 15] has not, to our knowledge, used a modern 063
vision-language model to predict the COPUS code set. 064

Re-interpreted as a vision-language task, COPUS has 065
four properties that current video-LM benchmarks largely 066
neglect. It provides dense, long-horizon structure (a 60 067
min lecture contains 30 intervals \times 24 codes = 720 pre- 068
dictions per video); an externally validated vocabulary stan- 069
dard across institutions and peer-reviewed [11]; a published 070
inter-rater reliability literature reporting aggregate Cohen’s 071
 κ of 0.79–0.87 across trained-observer pairs [11], giving an 072
empirical target for human-level agreement on the protocol; 073
and fine-grained, multi-modal distinctions (*e.g.* Clicker- 074
Group vs Worksheet-Group depends on small objects in the 075
scene) that are only partially observable from vision, moti- 076
vating audio fusion and cross-modal alignment work. 077

We propose COPUS as a multimodal foundation-model 078

Table 1. The 24 COPUS codes. Each interval receives a binary annotation per code.

Code	Who	Description	Code	Who	Description
L	Stu	Listening	Lec	Ins	Lecturing
Ind	Stu	Individual think.	RtW	Ins	Real-time writing
CG	Stu	Clicker group	FUp	Ins	Follow-up
WG	Stu	Worksheet group	PQ	Ins	Posing question
OG	Stu	Other group	CQ	Ins	Clicker question
AnQ	Stu	Answering Q.	AnQ	Ins	Answering Q.
SQ	Stu	Asks question	MG	Ins	Moving/guiding
WC	Stu	Whole-class disc.	1o1	Ins	One-on-one
Prd	Stu	Prediction	D/V	Ins	Demo/video
SP	Stu	Presentation	Adm	Ins	Administration
TQ	Stu	Test/quiz	W	Ins	Waiting
W	Stu	Waiting			
O	Stu	Other			

079 benchmark and characterise what distinguishes it from existing
080 video-LM evaluations; provide VISTA (Vision Instrument for STEM Teaching Activity), a MiniCPM-V-
081 4.5 [18] baseline with a structured multi-label prompt and max-pool aggregator matching the protocol’s any-
082 occurrence semantics; release evaluation tooling and the prompt vocabulary; and discuss failure modes that COPUS
083 isolates more cleanly than current video-LM benchmarks.
084
085
086

087 2. The COPUS Benchmark

088 Let \mathcal{A} denote the 24 COPUS codes (Table 1) and let $I_k = [120k, 120(k+1))$ seconds denote the k -th 2-minute interval
089 of a lecture of duration T seconds, $K = \lceil T/120 \rceil$. A system must produce, for each interval k and each $a \in \mathcal{A}$,
090 a binary label $y_{a,k} \in \{0, 1\}$. The reference matrix $Y \in \{0, 1\}^{|\mathcal{A}| \times K}$ is the consensus of five human evaluators, each
091 independently coding every lecture following the standard COPUS protocol [11]. Evaluation uses per-code accuracy
092 for non-rare codes, recall for rare codes (where accuracy is dominated by true negatives), and restricted macro accuracy
093 averaged over codes that occur in the evaluation set.
094
095
096
097
098

099 Our evaluation corpus consists of 13 human-annotated university chemistry lectures, recorded with a single fixed-
100 angle camera covering both instructor and students. Each lecture was independently coded by five trained human
101 evaluators using the COPUS spreadsheet format, and the resulting five evaluation matrices were merged into a single
102 consensus matrix that we treat as the ground truth Y . Our comparison tooling reads those matrices directly, so there is
103 no re-annotation step and no ground-truth transformation. For this paper we report results on three full lectures whose
104 code-frequency values span the full spectrum of rarity. We use the published COPUS inter-rater reliability literature as
105 the human noise floor for our corpus rather than computing Fleiss’ κ directly on our 5-rater matrices; the implications
106 of this choice are discussed in §5.
107
108
109
110
111
112
113

Four properties of COPUS surface failure modes that current short-form video-LM benchmarks do not isolate: multi-label co-occurrence (intervals can contain several simultaneous codes, yielding higher label density than Charades-style activity sets [8, 10]); fine-grained spatial reasoning; interval–event granularity mismatch, where events much shorter than the 2-minute interval must still be counted as present (encouraging localisation-style models [20]); and audio-partial observability, since question and discussion codes are defined partly by spoken intent rather than visible behaviour. Across our corpus, the codes LEC, L, RTW occur in >80% of intervals while ADM, SP, TQ occur in <5%. These are rare but important events that are statistically marginal in any naturally collected corpus.

Published reliability as context. Smith et al. [11] report aggregate Cohen’s κ of 0.79–0.87 across trained-observer pairs (averaged over all codes). Per-code reliability is reported as Jaccard similarity rather than κ because constant codes (always-absent or always-present) break κ ’s chance correction. Several rare codes (SP, PRD, TQ) did not occur in their validation corpus, so no published per-code statistic exists for them.

3. Baseline Pipeline

We provide a strong baseline (VISTA) whose modular structure makes three design decisions easy to vary independently: temporal granularity, prompt structure, and aggregation operator. The pipeline (Fig. 1) has three stages: 3FPS preprocessing, sliding-window vision-language model (VLM) inference with a multi-label prompt, and max-pool aggregation onto the 2-minute COPUS grid.

COPUS is defined over 2-minute intervals, but the events that determine a code’s presence (e.g. a 7 second RTW) may be much shorter. In the spirit of sparse temporal segment sampling [14], we use a dense sliding window $W_i = [is, is + \Delta t]$, $\Delta t=10$ s, $s=5$ s, giving 50% overlap and 719 windows per 60-minute lecture. Any behaviour lasting ≥ 5 s appears in multiple VLM predictions and the aggregation step sees redundancy.

One naive approach is to query each of the 24 codes separately, however this costs $24 \times$ more compute, is order-sensitive, and ignores joint structure. We instead give a single prompt per window instructing the VLM to generate observations for all 24 codes (Listing 1). On development data, the single-pass prompt produced more code coverage per window without any observable loss of precision, at 1/24 the inference cost.

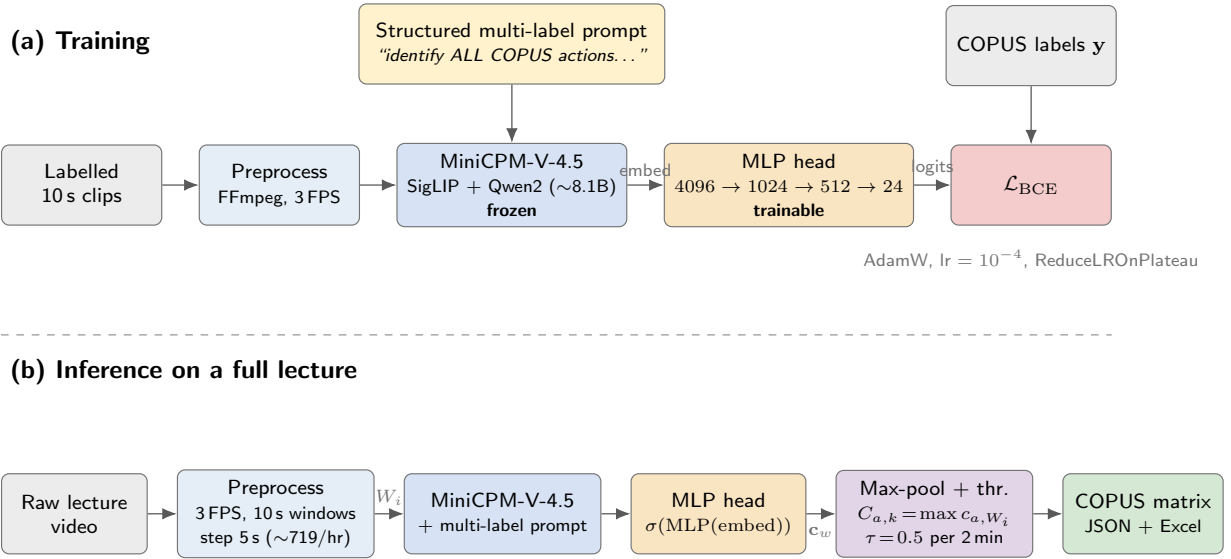


Figure 1. The VISTA baseline. (a) Training: a 3-layer multi-layer perceptron (MLP) head is trained on top of the frozen MiniCPM-V-4.5 backbone with binary cross-entropy (BCE) loss against multi-label clip annotations. The multi-label prompt is identical at training and inference. (b) Inference: a full lecture is preprocessed to 3 FPS and split into 10 s windows; per-window 24-dimensional confidences are max-pooled across windows overlapping each 2-minute COPUS interval and thresholded at $\tau = 0.5$.

Listing 1. Structured multi-label prompt.

```
Analyze this classroom video and identify ALL
COPUS actions that are currently occurring.

COPUS Actions to identify:
{copus_actions_list}

Instructions:
1. Watch the video segment carefully
2. Identify ALL actions that are happening (
   there can be multiple simultaneous actions)
3. For each action you identify, provide:
- The exact action code from the list above
- Your confidence level (high, medium, or low)
- A brief justification

Format your response as:
DETECTED ACTIONS:
[action_code_1]: [confidence_level] - [brief
justification]
[action_code_2]: [confidence_level] - [brief
justification]
...
```

The VLM’s response is converted to a 24-dim confidence vector $\mathbf{c}_w \in [0, 1]^{24}$ either by parsing the “DETECTED ACTIONS” block or by matching wording against a keyword dictionary based on previous model responses.

A 3-layer MLP ($4096 \rightarrow 1024 \rightarrow 512 \rightarrow 24$) is trained on top of the frozen VLM with BCE loss on a supervised clip corpus drawn from the 10 lectures held out of evaluation, supervised by the same 5-rater consensus. Training uses AdamW [6], base learning rate (LR) 10^{-4} with

ReduceLROnPlateau, weight decay 10^{-2} , and gradient clipping at 1.0. Because rare codes appear in $<5\%$ of intervals, a model trained on a uniform sample of intervals minimises BCE by predicting negative on every rare class. We address this by sampling the training clips for diversity: the training set deliberately mixes intervals with no codes, intervals with a single code, and intervals with multiple co-occurring codes, so the model has to learn both presence/absence discrimination and the joint structure of simultaneous behaviours. Full 24-way multi-label BCE supervision is applied to every clip, so the model is penalised for spurious positives on any of the 24 codes regardless of which clip type it sees.

Per-window confidences are grouped into the 2-minute COPUS grid by max-pooling over overlapping windows: $C_{a,k} = \max_{W_i \cap I_k \neq \emptyset} c_{a,W_i}$ and $\hat{y}_{a,k} = \mathbb{1}[C_{a,k} \geq 0.5]$. The choice of max over mean matches the protocol’s requirements: COPUS codes a behaviour as present if it occurs at any point in the interval, so mean pooling would penalise shorter actions/codes.

4. Results

Inference uses a single NVIDIA L40S. Thus, a 60-minute lecture produces 719 sliding windows at ~ 32 GPU-hours, with the dense windowing and 500-token structured-generation prompt deliberately expensive per window so that brief behaviours are not missed in long, high-FPS lecture video. We compare two configurations of VISTA:

161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183

185
186
187
188

189
190
191
192
193

194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213

214

215
216
217
218
219
220

Table 2. Per-code results across 3 lectures against the 5-rater consensus. Accuracy is reported for high- and mid-frequency codes, and recall is reported for rare codes. Restricted macro accuracy is averaged over the $m=10$ COPUS codes with ≥ 1 positive in the evaluation set. Codes absent from the corpus are excluded because predicting all-negative on an absent code trivially yields 100% accuracy and would inflate the macro.

Code	Ours (FT)	Ours (ZS)
<i>High-frequency (accuracy, %)</i>		
Lec	85.3	74.4
L	89.7	78.5
RtW	53.3	44.7
MG	100.0	91.9
<i>Mid-frequency (accuracy, %)</i>		
PQ	82.0	82.0
AnQ (Ins.)	76.0	74.0
CG / WG	78.0 / 72.0	79.0 / 75.0
<i>Rare (recall, %)</i>		
SP	80.0	70.0
Prd	70.0	65.0
SQ	70.0	65.0
Restricted macro acc. (m codes, ≥ 1 pos.)	80.1	74.9

221 a fine-tuned (FT) variant with the trained MLP head over
 222 MiniCPM-V-4.5 hidden states, and a zero-shot (ZS) variant
 223 using the keyword-parsed VLM output directly. MiniCPM-
 224 V-4.5 was chosen over other open-weights VLMs (Qwen2-
 225 VL [17], VideoLLaMA 2 [1], LLaVA-Video [21]) be-
 226 cause its SigLIP-style [19] visual encoder runs on-prem
 227 on consumer hardware, which is necessary for our access-
 228 restricted classroom video. Per-code accuracy is computed
 229 against the five-rater human consensus matrix. For rare
 230 codes (SP, PRD, SQ) we report recall rather than accuracy:
 231 because positive labels are sparse, a model that simply pre-
 232 dicted negative achieves $\geq 95\%$ accuracy on these codes de-
 233 spite identifying none of the actual occurrences.

234 The fine-tuned model agrees with the human consensus
 235 on 2013 of 2160 binary predictions across 3 lectures, and
 236 restricted macro accuracy (over the $m=10$ codes with ≥ 1
 237 positive in the evaluation set) is 80.1%. For context, the
 238 published aggregate Cohen’s κ for COPUS observer pairs
 239 is 0.79–0.87 [11], so our model’s residual error reflects
 240 model limitations rather than approaching the human noise
 241 floor. Disagreements are largely on visually similar instruc-
 242 tor codes (RTW alone accounts for the largest single-code
 243 error and appears similar to LEC and instructor ANQ) and
 244 on rare codes, where the model misses 20–30% of true oc-
 245 currences. These are likely due to limited exposure during
 246 pretraining, even with the diversity-balanced training-clip
 247 sampling described in §3. The zero-shot variant is below
 248 the fine-tuned model by roughly 9 pp on the high-frequency
 249 codes, is essentially tied on the mid-frequency codes, and
 250 drops 5-10 pp of recall on rare codes.

Three main failure modes. Audio-partial codes (SQ, in-
 251 structor ANQ, CQ): the model detects visual correlates
 252 (*e.g.* hand raise) but cannot distinguish question from com-
 253 ment. Fine-grained group-work confusion (CG/WG/OG):
 254 the codes differ only in the small object students hold, so the
 255 model fails to distinguish and defaults to the generic OG—a
 256 data-quality (camera placement) failure as much as a model
 257 failure. Long-tail recall (SP, PRD, TQ): the supervised cor-
 258 pus contains few examples and the VLM has likely seen
 259 few classroom “predictions” during its pretraining, yielding
 260 the classic data-mixture problem for rare but pedagogically
 261 important events.
 262

5. Discussion 263

The agenda for video-language benchmarks is increasingly
 264 limited not by data volume but by the challenge of creating
 265 tasks whose predictions have relevance outside the multi-
 266 modal community. COPUS demonstrates a complementary
 267 idea: benchmarks that have already been validated by a do-
 268 main community bring with them published reliability, op-
 269 erational definitions, and a downstream use case.
 270

Limitations. Several caveats apply. The evaluation cor-
 271 pus is narrow: three chemistry lectures from one institu-
 272 tion with one fixed-angle camera, so generalisation across
 273 STEM disciplines and recording setups is future work. We
 274 use the published COPUS reliability literature [11, 12] as
 275 the human noise floor rather than computing Fleiss’ κ on
 276 our 5-rater panel; because our evaluators followed the same
 277 protocol this is a reasonable proxy, but a direct κ on our
 278 matrices is the natural follow-up. The label distribution is a
 279 significant data challenge in its own right: with some codes
 280 appearing in less than 5% of intervals, an unconstrained
 281 training objective can be trivialised by predicting negative
 282 everywhere, so recall must be used to evaluate performance.
 283 Our diversity-balanced training-clip sampling mitigates this
 284 but does not eliminate it; constructing better-balanced train-
 285 ing data is a clear direction. We have not yet compared
 286 to a non-VLM baseline such as VideoMAE [13] nor per-
 287 formed extensive audio fusion. We have institutional access
 288 to the classroom video but can not release it publicly. The
 289 open-weights backbone [18] lets the evaluation run on local
 290 hardware, and we release the pipeline, prompts and eval-
 291 uation tooling but not the raw video. The system should
 292 not be used to evaluate individual instructors: published
 293 κ on several codes is too low for fair single-person deci-
 294 sions. COPUS-style imported benchmarks offer a low-cost
 295 path to reliability-calibrated multimodal evaluation, and the
 296 same template should apply to other domains where vali-
 297 dated coding schemes already exist.
 298

299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355**References**

- [1] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. VideoLLaMA 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. <https://arxiv.org/abs/2406.07476>. 4
- [2] Patrick J. Donnelly, Nathaniel Blanchard, Borhan Samei, Andrew M. Olney, Xiaoyi Sun, Brooke Ward, Sean Kelly, Martin Nystrand, and Sidney K. D’Mello. Multi-sensor modeling of teacher instructional segments in live classrooms. In *Proc. 18th ACM International Conference on Multimodal Interaction (ICMI)*, 2016. 1
- [3] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-MME: The first-ever comprehensive evaluation benchmark of multi-modal LLMs in video analysis. *arXiv preprint arXiv:2405.21075*, 2024. <https://arxiv.org/abs/2405.21075>. 1
- [4] Matthew T. Hora and Joseph J. Ferrare. The teaching dimensions observation protocol (TDOP) 2.0 user’s guide. Technical report, Wisconsin Center for Education Research, University of Wisconsin–Madison, 2014. <https://tdop.wceruw.org/TDOP-2.1-Users-Guide.pdf>. 1
- [5] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, Limin Wang, and Yu Qiao. MVBenCh: A comprehensive multi-modal video understanding benchmark. In *CVPR*, 2024. arXiv:2311.17005. 1
- [6] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. arXiv:1711.05101. 3
- [7] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. EgoSchema: A diagnostic benchmark for very long-form video language understanding. In *NeurIPS*, 2023. arXiv:2308.09126. 1
- [8] A. J. Piergiovanni and Michael S. Ryoo. Learning latent super-events to detect multiple activities in videos. In *CVPR*, 2018. https://openaccess.thecvf.com/content_cvpr_2018/html/Piergiovanni_Learning_Latent_Super-Events_CVPR_2018_paper.html. 2
- [9] Daiyo Sawada, Michael D. Piburn, Eugene Judson, Jeff Turley, Kathleen Falconer, Russell Benford, and Irene Bloom. Measuring reform practices in science and mathematics classrooms: The reformed teaching observation protocol. *School Science and Mathematics*, 102(6):245–253, 2002. 1
- [10] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, 2016. arXiv:1604.01753. 2
- [11] Michelle K. Smith, Francis H. M. Jones, Sarah L. Gilbert, and Carl E. Wieman. The classroom observation protocol for undergraduate STEM (COPUS): A new instrument to characterize university STEM classroom practices. *CBE—Life Sciences Education*, 12(4):618–627, 2013. 1, 2, 4
- [12] Marilyne Stains, J. Harshman, M. K. Barker, S. V. Chasteen, R. Cole, S. E. DeChenne-Peters, M. K. Eagan, J. M. Esson, J. K. Knight, F. A. Laski, et al. Anatomy of STEM teaching in North American universities. *Science*, 359(6383):1468–1470, 2018. 4
- [13] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *NeurIPS*, 2022. arXiv:2203.12602. 4
- [14] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. arXiv:1608.00859. 2
- [15] Zuowei Wang, Xingyu Pan, Kevin F. Miller, and Kai S. Cortina. Automatic classification of activities in classroom discourse. *Computers & Education*, 78:115–123, 2014. 1
- [16] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. LongVideoBench: A benchmark for long-context interleaved video-language understanding. In *NeurIPS*, 2024. arXiv:2407.15754. 1
- [17] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. <https://arxiv.org/abs/2407.10671>. 4
- [18] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. MiniCPM-V: A GPT-4V level MLLM on your phone. *arXiv preprint arXiv:2408.01800*, 2024. <https://arxiv.org/abs/2408.01800>. 2, 4
- [19] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. arXiv:2303.15343. 4
- [20] Chen-Lin Zhang, Jianxin Wu, and Yin Li. ActionFormer: Localizing moments of actions with transformers. In *ECCV*, 2022. arXiv:2202.07925. 2
- [21] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. LLaVA-Video: Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. <https://arxiv.org/abs/2410.02713>. 4