# VL-JEPA: JOINT EMBEDDING PREDICTIVE ARCHITECTURE FOR VISION-LANGUAGE

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

018

019

021

024

025

026027028

029

031

032

033

034

037

040

041

042

043

044

046

047

048

050 051

052

Paper under double-blind review

### **ABSTRACT**

We introduce VL-JEPA, a vision-language model built on a Joint Embedding Predictive Architecture (JEPA). Instead of autoregressively generating tokens as in classical VLMs, VL-JEPA predicts continuous embeddings of the target texts. By learning in an abstract representation space, the model can focus on task-relevant semantics while abstracting away surface-level linguistic variability. In a strictly controlled comparison against standard token-space VLM training with the same vision encoder and training data, VL-JEPA achieves stronger performance while having 50% fewer trainable parameters. At inference time, a lightweight text decoder is invoked only when needed to translate VL-JEPA predicted embeddings into text. We show that VL-JEPA natively supports selective decoding that can reduce the number of decoding operations by  $\sim 2.85 \times$  while maintaining similar performance compared to dense non-adaptive uniform decoding. Beyond generation, the embedding-space formulation naturally supports open-vocabulary classification and retrieval without any architecture modification. VL-JEPA achieves leading SoTA zero-shot results on diverse real-world video-language understanding tasks on COIN, CrossTask, EgoExo4D, SSv2, and WORLDPREDICTION-WM, substantially outperforming larger generative VLMs trained with more data.

#### 1 Introduction

One of the most important aspects of advanced machine intelligence is the ability to understand the physical world that surrounds us. This ability enables AI systems to learn, reason, plan and act in the real world in order to assist humans (LeCun, 2022). Intelligent systems that need to act in the real world includes robots and wearable devices (Fung et al., 2025). Machine learning tasks that make up for this ability include captioning, retrieval, visual question answering, action tracking, reasoning and planning etc (Bordes et al., 2024; Chen et al., 2025b). Other systems requirements for real-world applications include real-time response with low latency and models with lower inference cost.

Currently, the common approach to achieve these tasks is to train large token-generative Vision Language Models (VLMs) (Liu et al., 2023; Dai et al., 2023; Alayrac et al., 2022; Chen et al., 2024b; Cho et al., 2025; Chen et al., 2022), which takes visual input  $X_V$ , textual query  $X_Q$  to generate desired textual response Y autoregressively in token space, i.e.,  $(X_V, X_Q) \mapsto Y$ . This is straightforward but inadequate for two main reasons. First, VLMs are expensive to develop, because they are trained to generate responses Y to queries by capturing both task-relevant semantics with task-irrelevant surface linguistic features such as words choice, style or paraphrasing. During training, VLMs must model both aspects, which results in unnecessary computing effort spent producing diverse token sequences that ultimately do not impact the correctness of the output. Second, real-time tasks involving live streaming video (e.g., live action tracking) require sparse and selective decoding (e.g.,, emitting a description only when a new event occurs) (Zhou et al., 2024). However, VLMs rely on autoregressive token-by-token decoding, which must be completed before revealing the underlying semantics of Y. This process introduces unnecessary latency and hampers the ability to update semantics dynamically in real time.

This paper introduces a Joint Embedding Predictive Architecture for Vision-Language (VL-JEPA), turning expensive learning of data-space token generation into more efficient latent-space semantic prediction. As illustrated in Fig. 1, the model employs **x-encoder** to map vision inputs  $X_V$  into embedding  $S_V$ , a **y-encoder** to map the textual target Y into an embedding  $S_Y$ , and a **predictor** 

that learns the mapping  $(S_V, X_Q) \mapsto S_Y$  where  $X_Q$  is a textual query (i.e., the prompt). The training objective is defined in the embedding space  $\mathcal{L}_{\text{VL-JEPA}} = D(\hat{S}_Y, S_Y)$  instead of the data space  $\mathcal{L}_{\text{VLM}} = D(\hat{Y}, Y)$ . During inference, a **y-decoder** reads out the predicted embedding  $\hat{S}_Y$  to text space  $\hat{Y}$  when needed.

Thanks to its **non-generative** nature, VL-JEPA is not forced to reconstruct every surface detail of Y in the token space. Instead, it only needs to predict the abstract representation  $S_Y$  in the embedding space. In the raw one-hot token space, different plausible Y outputs for the same input often appear nearly orthogonal if they don't share overlapping tokens. However, in the embedding space, these diverse targets can be mapped to nearby points that share similar semantics. This simplifies the target distribution thus makes the learning process more efficient. In addition, unlike VLMs, this approach eliminates the need for learning language generation with a heavy decoder during training, resulting in significant efficiency gains.

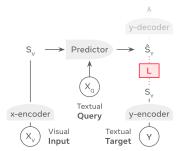


Figure 1: VL-JEPA.

Thanks to its **non-autoregressive** nature, VL-JEPA can produce continuous streams of target semantic embeddings within sliding windows with minimal latency as it only require a single forward pass without autoregressive decoding. This is particularly advantageous for real-time online applications such as live action tracking, scene recognition, or planning, where the embedding stream can be selectively decoded by a lightweight y-decoder, enabling efficient and prompt updates.

In this work, we empirically validate the advantages of VL-JEPA. First, we conduct a strictly controlled comparison against classical token-generative VLM (Liu et al., 2023; Cho et al., 2025): both setups use the same vision encoder, spatial resolution, frame rate, training data, batch size, and number of iterations, etc., with the *only* difference being the objective in token space or embedding space. Under this matched training condition, VL-JEPA delivers consistently higher performance on zero-shot captioning and classification while using roughly half the trainable parameters, indicating that embedding-space supervision improves learning efficiency.

Beyond the training phase, VL-JEPA also delivers substantial inference-time efficiency improvement through *selective decoding*, where decoding happens only due to significant change in the predicted embedding stream. Empirically, this strategy reduces the number of decoding operations by  $\sim 2.85 \times$  while preserving overall output quality, as measured by average CIDEr scores against human annotations.

VL-JEPA achieves leading *zero-shot* performance across a diverse set of video-language benchmarks, including step recognition on COIN (Tang et al., 2019) and CrossTask (Zhukov et al., 2019), fine-grained keystep recognition on EgoExo4D (Grauman et al., 2024), and video classification on SSv2 (Goyal et al., 2017), SSv2-Events (Bagad et al., 2023), and SSv2-Temporal (Sevilla-Lara et al., 2021). It further demonstrates a non-trivial zero-shot text-to-video retrieval performance on YouCook2 (Zhou et al., 2018), achieving 17.7 recall@1, while not training with contrastive objective. All these results are obtained with a single VL-JEPA model, without any architectural modifications or task-specific heads. Beyond zero-shot evaluations, we further show that short finetuning of VL-JEPA yields state-of-the-art results on COIN step recognition, reaching 72.84% accuracy. The code and model will be open-sourced. In summary, the contributions of this paper are as follows:

- We introduce VL-JEPA, the first non-generative model that can perform general-domain vision-language tasks in real-time, built on a joint embedding predictive architecture. It achieves SoTA zero-shot scores on wide range of video-language understanding tasks.
- We demonstrate in controlled experiments that VL-JEPA, trained with latent space embedding prediction, outperforms VLMs that rely on data space token prediction.
- We show that VL-JEPA delivers significant efficiency gains over VLMs for online video streaming applications, thanks to its non-autoregressive design and native support for selective decoding.
- We demonstrate the scalability of VL-JEPA with consistent improvement when scaling the number of model parameters and dataset size.

## 2 RELATED WORKS

**JEPA Models.** JEPA model learns by predicting the representation of a target input Y from the representation of a context input X. Early instantiations include I-JEPA for image encoding (Assran et al., 2023) and V-JEPA for video encoding (Bardes et al., 2023), which demonstrated the effectiveness of this objective over pixel reconstruction approach in their respective modality. Recent JEPA work falls into two categories. One category of work emphasizes better unimodal representation learning (Assran et al., 2023; Bardes et al., 2023; Fei et al., 2023) or cross-modal alignment (Lei et al., 2025; Jose et al., 2025). The other direction targets world modeling, where pretrained encoders are frozen and action-conditioned predictors are trained for conditional prediction of state representations (Zhou et al., 2025; Baldassarre et al., 2025; Assran et al., 2025). This has shown good results but remains limited to narrow domains like mazes or robotic pick-and-place. Our proposed VL-JEPA is the first designed for general-purpose vision—language tasks. It performs conditional latent prediction over vision and text, and preserves efficiency while enabling flexible, multitask architecture.

Vision Language Models. Existing vision-language models largely fall into two families: (1) CLIP-style models with a non-predictive joint-embedding architecture (JEA) (Radford et al., 2021; Zhai et al., 2023; Bolya et al., 2025; Liu et al., 2024; Chen et al., 2023) encode images and texts independently into a common latent space,  $X_V \mapsto S_V$  and  $Y \mapsto S_Y$ . By minimizing  $\mathcal{L}_{\text{CLIP}} = D(S_V, S_Y)$  with a contrastive loss (e.g., InfoNCE), CLIP learns aligned representations that support zero-shot classification and vision-language retrieval; (2) Generative VLMs (Liu et al., 2023; Chen et al., 2022; Dai et al., 2023; Alayrac et al., 2022; Chen et al., 2024b; Cho et al., 2025; Beyer et al., 2024) connect a vision encoder (Radford et al., 2021; Fini et al., 2025) with a language model (e.g., LLM). They are typically trained with  $\mathcal{L}_{\text{VLM}} = D(\hat{Y}, Y)$ , i.e., next token prediction with crossentropy loss, and can learn to handle various vision-text-to-text generation tasks such as visual question answering (VQA).

Our proposed VL-JEPA integrates the architectural advantages and task coverage of both CLIPs and VLMs (Table 1). Since VL-JEPA learns in embedding space, it can leverage web-scale noisy image—text pairs (Jia et al., 2021), yielding strong open-domain features. On the other hand, VL-JEPA supports conditional generation tasks with a readout text decoder. Meanwhile, compared to generative VLMs



Table 1: Task coverage comparison.

that optimize directly in data space, VL-JEPA is more efficient at learning in the latent space. In addition, it is also more efficient for online inference, as it allows naturally selective decoding.

## 3 METHODOLOGY

**Model Architecture.** We propose **VL-JEPA** (Fig. 2), which instantiates the joint embedding predictive architecture (JEPA) for vision–language learning. VL-JEPA learns triplets  $\langle X_V, X_Q, Y \rangle$ , where  $X_V$  denotes the visual input (a single image or a sequence of video frames), and  $(X_Q, Y)$  are the textual query and target (e.g., an instruction and its answer). The model comprises of four modules:

- 1. **Vision encoder** compresses high-volume visual inputs to compact representations,  $X_V \mapsto S_V$ , where  $S_V$  is a sequence of continuous vectors analogous to "visual tokens" in VLMs. VL-JEPA is agnostic to the specific encoder: one can use image encoders such as CLIP (Radford et al., 2021), Perception Encoder (Bolya et al., 2025), DINOv2/v3 (Oquab et al., 2023; Siméoni et al., 2025), or video encoders such as V-JEPA 2 (Assran et al., 2025).
- 2. **Text encoder** embeds the target into a latent space,  $Y \mapsto S_Y$ , serving as the y-encoder of JEPA. It also encodes the query text to embeddings, which is fed into the predictor along with visual tokens. In this work, the text encoder is kept frozen to avoid representation collapse.
- 3. **Predictor** is the core component of VL-JEPA. It learns the mapping  $(S_V, X_Q) \mapsto S_Y$ . Visual and query embeddings are first projected into a common dimension and then passed through multiple Transformer layers. After pooling (e.g., average pooling or selecting the query token), the output is projected into the  $S_Y$  space.

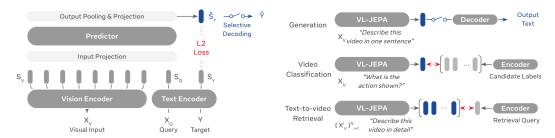


Figure 2: **Left**: VL-JEPA learns to predict target embedding  $S_Y$  instead of reconstruct raw target Y in token space as in classical VLMs. VL-JEPA can do selective decoding during inference, based on the predicted embedding stream. **Right**: In addition to vision-text-to-text tasks such as captioning, the VL-JEPA's embedding space allows to also handle open-vocabulary classification and text-to-video retrieval tasks using a single model architecture.

4. **Text decoder** is not involved during training. At inference time, it reads out the predicted embedding as human-readable text,  $\hat{S}_Y \mapsto \hat{Y}$ .

**Training Objective.** With the target embedding  $S_Y = \text{y-encoder}(Y)$  and the prediction  $\hat{S}_Y$ , the training loss is defined as the L2 distance between them:  $\mathcal{L}_{\text{VL-JEPA}} = \|\hat{S}_Y - S_Y\|_2^2$ .

Compared to the token-space cross-entropy loss used by generative VLMs, doing regression in the embedding space benefits from a *simplified target distribution*. Specifically, many real-world prediction tasks are inherently ill-posed: for the same input X, there may exist multiple plausible targets Y that are all acceptable. For example, given the query "What will happen if I flip this light switch down?", both "the lamp is turned off" and "room will go dark" are valid answers. If VL-JEPA's y-encoder embeds them into nearby points (ideally yielding a compact unimodal distribution) the learning task will become much easier: the model no longer needs to fit multiple disjoint high-density regions in sparse token space, but only a single coherent mode in a continuous embedding space. Additionally, this loss does not require running inference and backpropagation in the LLM decoder during training, yielding considerable efficiency gains.

However, embedding space loss raises the risk of representation collapse, where the y-encoder could over-simplify  $S_Y$  to constant vectors, allowing the predictor to learn only a trivial mapping. Methods to address this issue involves adding regularization term (Bardes et al., 2021), employing exponential moving average (EMA) to y-encoder (Assran et al., 2025), or freezing the y-encoder (Zhou et al., 2025). In this work, we adopt the last strategy for 1) simplicity, 2) training efficiency, and 3) the availability of text encoders that effectively capture task-relevant semantics while filtering out surface-level linguistic details.

Multi-tasking with a Single Architecture. VL-JEPA performs query-conditioned prediction in embedding space, allowing a single architecture to support diverse task families (Fig. 2). For vision-to-text generation tasks such as captioning, the query  $X_Q$  is a captioning prompt and the predictor learns to predict the embedding of the target caption, which is then decoded into text. VL-JEPA also supports CLIP-style open-vocabulary classification and text-to-video retrieval. In classification, candidate label texts are encoded into embeddings and compared with  $\hat{S}_Y$  to select the nearest match. In retrieval, candidate videos are mapped to  $\hat{S}_Y$  with a captioning prompt and ranked by similarity to the encoded query.

Selective Decoding for Streaming Video Applications. Real-world video applications often require online streaming inference, such as tracking user actions in smart glasses for procedural assistance (Chen et al., 2024c), monitoring world states for online planning, navigation and robotics (Shukor et al., 2025; Black et al., 2025; Song et al., 2025). A central challenge is balancing two competing needs: on the one hand, the model must continuously update semantics as new frames arrive; on the other hand, computational efficiency and latency are critical. Existing VLMs typically rely on explicit memory mechanisms (Zhou et al., 2024; Qian et al., 2024) to decide when to decode

or complex KV-cache optimizations (Di et al., 2025) for efficiency, since autoregressive language models are expensive to run continuously.

VL-JEPA, in contrast, natively supports selective decoding. Since it predicts a semantic answer embedding non-autoregressively, the model provides a continuous semantic stream of  $\hat{S}_Y$  that can be monitored in real time. This stream can be stabilized with simple smoothing (e.g., average pooling) and decoded only when a significant semantic shift is detected, such as when the local window variance exceeds a threshold. In this way, VL-JEPA maintains always-on semantic monitoring while avoiding unnecessary decoding, achieving both responsiveness and efficiency.

## 4 EXPERIMENTS

We begin by outlining the implementation details of VL-JEPA in §4.1. In §4.2, we demonstrate the advantage of embedding prediction by comparing it with a token-predictive VLM baseline under a strictly controlled setting. In §4.3, we evaluate the effectiveness of VL-JEPA's selective decoding, and show that it reduces decoding cost while maintaining the performance. Next, we benchmark VL-JEPA against state-of-the-art models across a range of downstream tasks, including zero-shot and finetuning video understanding (§4.4, §4.5), and WorldPrediction (§4.6). Finally, we present ablation studies in §4.7 and scalability analysis in §4.8. Additional experimental details are deferred to the appendix.

#### 4.1 IMPLEMENTATION OF VL-JEPA

**Data.** We include <u>PLM-IMAGE-AUTO</u>, which is one of the core components of PLM's training data (Cho et al., 2025). It provides detailed captions generated by a synthetic engine powered by Llama-3.2-90B-Vision. We include its SA-1B split (9.35M) and OpenImages split (1.64M). <u>DATACOMP</u> (Gadre et al., 2023) is a large-scale web image—text collection commonly used for vision-language pretraining; we use 10.1M re-captioned samples from (Li et al., 2024). <u>PIXMO-CAP</u> (Deitke et al., 2025) is a high-quality human-annotated dataset frequently used in recent VLMs, from which we use 0.59M detailed captions. We include video captioning data from <u>PLM-VIDEO-AUTO</u> (2.1M from YT-1B and 0.17M from Ego4D) (Cho et al., 2025). We also include atomic action descriptions from <u>Ego4D</u> (Grauman et al., 2022) (3.70M) and <u>EgoExo4D</u> (Grauman et al., 2024) (0.95M). These fine-grained action descriptions are commonly used to train egocentric vision language models, *e.g.*, EgoVLP (Lin et al., 2022a), LaViLa (Zhao et al., 2023). <u>ACTION100M</u>. An internal dataset annotated on HowTo100M following the methodology used in VLWM (Chen et al., 2025b). It contains 100M automatically annotated segments across 0.71M instructional videos. For each segment, we uniformly sample one annotation from {brief caption, detailed caption, brief action, detailed action}.

**Model.** <u>Vision Encoder.</u> Unless otherwise specified, we use a frozen V-JEPA 2 ViT-L (Assran et al., 2025) with 304M parameters. Each video input is uniformly sampled into 64 frames at 256<sup>2</sup> resolution. For image inputs, the same image is duplicated 64 times to match the input shape. The encoder outputs 8192 visual tokens, each of dimension 1024. <u>Text Encoder & Decoder.</u> We use a pretrained encoder–decoder language model capable of both encoding and decoding nearly identical text representations. For the experiments in this paper, we adopt a variant of the SONAR models (Duquenne et al., 2023), trained with a sequence-to-sequence NLL objective and an additional InfoNCE loss. The model can be easily replaced with other off-the-shelf alternatives. The encoder outputs two 1024-dimensional embeddings—one for the query and one for the target—which remain frozen during pretraining. The decoder is used only at inference time.

<u>Predictor.</u> The predictor is initialized with the last 8 Transformer layers of Llama-3.2-1B, resulting in 490M trainable parameters. We disable the causal attention mask so that both vision and query embeddings can be jointly attended. Linear projections connect the predictor with the vision and text encoders, and average pooling is applied to obtain the predicted target embedding.

**Training Setup.** VL-JEPA is trained on 16 nodes with  $8\times NVIDIA$  H200 GPUs each, using a global batch size of 2048 and bf16 precision. We adopt a constant learning rate of  $1\times 10^{-4}$  to facilitate stable resumption and extended training, though better scheduling like cosine may further improve performance. Models are trained for 30k iterations, corresponding to 61.4M seen samples.

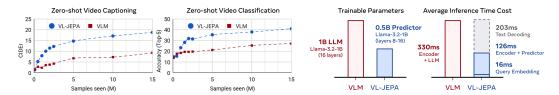


Figure 3: Comparison of embedding prediction (VL-JEPA) and token prediction (VLM). We conduct a fair comparison of under strictly aligned training settings (encoder, data, batchsize, etc.). Left: Zero-shot video captioning CIDEr score averaged over 3 datasets and zero-shot classification accuracy (top-5) averaged over 3 benchmarks. Right: Comparing the trainable parameters and average inference time cost.

### 4.2 EMBEDDING PREDICTION VS. TOKEN PREDICTION: A CONTROLLED COMPARISON

**Evaluation Setup.** In this section, we compare VL-JEPA to a token-generative VLM baseline under a strictly aligned training conditions. Both models use the same Perception Encoder (Bolya et al., 2025) (frozen ViT-L-14 with 336² resolution, no tiling, 16 frames per video) for vision inputs. We use the same training iterations with the same effective batch size of 128, same learning rate scheduler on the same pretraining data mixture described above (§4.1). The only difference is the prediction task: VL-JEPA predicts target embeddings (Duquenne et al., 2023) using a 0.5B predictor, whereas the VLM baseline performs next-token prediction with cross-entropy using a 1B LLM. For VLM, we use the standard training recipe and codebase of PerceptionLM (Cho et al., 2025), aligning frozen vision encoder and text-only LLM Llama-3.2-1B. For VL-JEPA, we initialize the predictor from the 8-16 layers of Llama-3.2-1B.

We evaluate both models at regular checkpoints throughout training spanning from 500K to 15M samples seen. At each checkpoint, we measure the performance on video captioning and video classification. For video captioning, we report CIDEr scores averaged across YouCook2 (Zhou et al., 2018), MSR-VTT (Xu et al., 2016) and PVD-Bench (Bolya et al., 2025). VL-JEPA decodes the predicted embeddings while VLM generates the tokens directly. For video classification, we report top-5 accuracy averaged across CrossTask-Step, CrossTask-Task (Zhukov et al., 2019) and EgoExo4D (Grauman et al., 2024). For VL-JEPA we choose the candidate with lowest cosine distance to the predicted embedding, while for VLM we pick the class with lowest perplexity.

**Results.** As shown in Fig. 3, both models yield comparable performance after 500K samples seen in both tasks, with respectively 1.23 and 1.35 CIDEr in video captioning and 14.9% and 14.0% top-5 accuracy for VL-JEPA and VLM. After a few iterations, we show that VL-JEPA's performance increase is much sharper compared to VLM, reaching 14.7 CIDEr and 35.3% top-5 accuracy after 5M samples seen. This gap remains constant as training scales at 15M samples with 14.8 CIDEr and 41.0% top-5 accuracy for VL-JEPA, while the VLM baseline yield respectively 7.1 CIDEr and 27.2% top-5 accuracy. This controlled comparison highlights the benefit of predicting embeddings rather than tokens, showing both higher sample efficiency and stronger absolute performance.

We compare the inference cost of the above VL-JEPA and the VLM by pre-loading 64 video frames into memory and repeatedly decoding text 100 times with the same prompt, measuring the average time per sample. As shown in Fig. 3 (right most), both models exhibit comparable latency when generating text. What differentiates our model from classical VLM is the decoupling between the prompt processing ("Query Embedding") and the video encoder ("Encoder + Predictor") from the text generation module ("Decoder"). This allows us to only use the first part of the model to perform retrieval and decode text only when needed (see Section 4.3 below), making our model more scalable for online video inference.

#### 4.3 EFFECTIVENESS OF SELECTIVE DECODING

**Evaluation Setup.** We evaluate the effectiveness of VL-JEPA's embedding-guided selective decoding on long-form video streams. To this end, we design a benchmark task where the goal is to recover a temporal sequence of annotations while minimizing the number of text decoding operations, which dominate inference cost. As shown in Fig. 4 (left), decoding is performed only

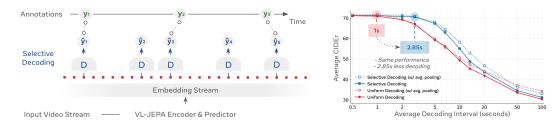


Figure 4: **Evaluation of selective decoding. Left:** We compare uniform sampling of decoding points at fixed intervals (red) and embedding-guided selective decoding (blue). Performance is measured by the average CIDEr score between each annotation y and its closest decoded output  $\hat{y}$ . **Right:** Results on EgoExo4D show that selective decoding achieves a Pareto improvement over uniform sampling: for the same performance level, it requires fewer decoding operations.

at selected points along the VL-JEPA embedding stream, yielding a sequence of N decoded outputs  $[(\hat{t}_1,\hat{y}_1),(\hat{t}_2,\hat{y}_2),\dots,(\hat{t}_N,\hat{y}_N)]$ . Each ground-truth annotation  $[(t_1,y_1),(t_2,y_2),\dots,(t_T,y_T)]$  is then aligned to its nearest decoded output in time (illustrated as  $\circ \cdots \circ$  in Fig. 4), and CIDEr is computed between matched pairs. We use the EgoExo4D (Grauman et al., 2024) validation set in procedural activity domains, which consists of 218 videos with an average duration of 6 minutes and about T=143 atomic action annotations per video.

As a baseline, we consider *uniform sampling*, where decoding points are placed at fixed intervals regardless of the underlying video content. Standard streaming VLMs are limited to this strategy, whereas VL-JEPA supports a more effective alternative: *adaptive selection* of decoding points guided by its predicted embeddings. We apply agglomerative clustering with temporal connectivity constraints (Murtagh & Contreras, 2012) to partition the embedding sequence into N segments of high intra-segment monosemanticity (Chen et al., 2024a), measured by variance (*i.e.*, Ward distance). The intuition is that within a semantically coherent segment, decoded outputs are highly similar, so decoding once per segment captures the essential information while greatly reducing overall decoding cost. The midpoint of each segment is then chosen as the decoding point, and decoding is performed either from the exact embedding or from the average-pooled embedding within the segment.

**Results.** As shown in Fig. 4 (right), we sweep the average decoding frequency from 2.0 Hz down to 0.01 Hz (*i.e.*, average intervals between consecutive decoding operations from 0.5s to 100s) by adjusting either the stride of uniform sampling or the number of clusters in adaptive selection. Across the entire range, adaptive selection consistently Pareto-dominates uniform sampling. In particular, selective decoding at 0.35 Hz (*i.e.*,  $\sim$ 2.85s interval) matches the performance of uniform decoding at 1 Hz, reducing decoding cost by  $\sim$ 2.85×. We further observe that average pooling provides consistent gains for both strategies, since it provides denoising and stabilization on embeddings prior feeding into the decoder.

#### 4.4 ZERO-SHOT VIDEO UNDERSTANDING

**Evaluation Setup.** We evaluate VL-JEPA following the CLIP-style evaluation protocol: candidate labels are embedded with the target encoder and matched against the predicted embeddings using cosine similarity. We assess VL-JEPA on a broad suite of video classification benchmarks. We evaluate on COIN (Tang et al., 2019) and CrossTask (Zhukov et al., 2019), two widely used dataset constructed from YouTube instructional videos. We include EgoExo4D fine-grained keystep recognition benchmark (Grauman et al., 2024), EPIC-KITCHENS-100 (EK-100) (Damen et al., 2022) action recognition benchmark, and Something-Something-v2 (SSv2) (Goyal et al., 2017) classification. In addition, we assess zero-shot classification on SSv2-Events (Bagad et al., 2023) and SSv2-Temporal (Sevilla-Lara et al., 2021).

**Results.** Table 2 reports zero-shot results. VL-JEPA surpasses existing baselines on COIN, CrossTask, EgoExo4D, and SSv2. On EK-100, it outperforms LaViLa (Zhao et al., 2023) and PE models and approaches the performance of GPT4Ego (Dai et al., 2024), which relies on a LaV-

COIN Step Recognition		CrossTask Step Recognition	on	EgoExo4D Keystep Recognition		
Model	Acc.	Model	Acc.	Model	Acc.	
Random Performance	0.1	Random Performance	1.0	Random Performance	0.4	
DistantSup. (Lin et al., 2022b)	10.2	PE-Core-L (Bolya et al., 2025)	38.4	PE-Core-L (Bolya et al., 2025)	10.9	
CLIP (Radford et al., 2021)	14.8	PE-Core-G (Bolya et al., 2025)	42.1	PE-Core-G (Bolya et al., 2025)	12.9	
ProcedureVRL (Zhong et al., 2023)	16.6	Qwen2.5VL 7B (Bai et al., 2025)	31.9	Qwen2.5VL 7B (Bai et al., 2025)	14.6	
VL-JEPA	16.7	VL-JEPA	49.0	VL-JEPA	22.2	

EPIC-KITCHENS-100 A	ction R	lecogni	tion	Something-Something-v2 Video Classification					
Model	Verb	Noun	Noun Action Model		SSv2	SSv2-Temporal	SSv2-Events		
Random Performance	1.3	0.5	0.1	Random Performance	0.6	5.6	2.0		
PE-Core-B (Bolya et al., 2025)	7.8	10.0	3.3	PE-Core-B (Bolya et al., 2025)	5.8	18.5	12.2		
PE-Core-L (Bolya et al., 2025)	10.7	17.3	6.0	PE-Core-L (Bolya et al., 2025)	9.3	31.9	16.2		
PE-Core-G (Bolya et al., 2025)	9.9	19.5	6.5	PE-Core-G (Bolya et al., 2025)	9.1	29.6	18.2		
Qwen2.5VL 3B (Bai et al., 2025)	18.7	11.8	5.7	Owen2.5VL 3B (Bai et al., 2025)	2.9	26.9	6.1		
Qwen2.5VL 7B (Bai et al., 2025)	15.0	15.8	5.7	Owen2.5VL 7B (Bai et al., 2025)	7.5	26.4	9.9		
Qwen2.5VL 32B (Bai et al., 2025)	18.3	20.6	8.6	Owen2.5VL 32B (Bai et al., 2025)	9.5	28.4	10.0		
LaViLa-B (Zhao et al., 2023)	_	_	16.3	VideoCLIP-B (Xu et al., 2021)	_	9.8	6.4		
LaViLa-L (Zhao et al., 2023)	_	-	23.8	VideoCon-L (Bansal et al., 2024)	_	15.2	11.4		
GPT4Ego-B (Dai et al., 2024)	_	_	28.9	VideoPrism-B (Zhao et al., 2024)	_	16.1	11.9		
GPT4Ego-L (Dai et al., 2024)	_	_	33.2	VideoPrism-g (Zhao et al., 2024)	_	18.6	15.7		
VL-JEPA	24.3	35.4	<u>25.0</u>	VL-JEPA	10.6	29.6	15.8		

Table 2: Zero-shot video understanding benchmark results.

iLa backbone combined with a heavy pipeline involving ChatGPT chain-of-thought prompting and SAM segmentation (Kirillov et al., 2023). In contrast, VL-JEPA achieves competitive accuracy with a much simpler and more efficient design.

#### 4.5 FINETUNING VL-JEPA

In addition zero-shot evaluation, we further assess the ability of VL-JEPA to adapt to downstream tasks through short finetuning (less than 10 epochs). We use the same L2 loss for fine-tuning. Table 3 shows results of full-shot step recognition on COIN (Tang et al., 2019) and CrossTask (Zhukov et al., 2019). VL-JEPA achieves state-of-the-art performance, reaching 77.2% on COIN and 86.9% on CrossTask, surpassing previous methods that rely on large VLMs.

Table 3: **Finetuning step recognition results.** 

Model	COIN	CrossTask
ClipBERT (Lei et al., 2021)	30.8	_
VideoCLIP (Xu et al., 2021)	51.2	60.1
TimeSformer (Bertasius et al., 2021)	54.6	60.9
DistantSup. (Lin et al., 2022b)	57.0	64.2
VideoTaskGraph (Ashutosh et al., 2023)	57.2	64.5
Paprika (Zhou et al., 2023)	51.0	63.5
VideoTF (Narasimhan et al., 2023)	56.5	_
ProcedureVRL (Zhong et al., 2023)	56.9	_
VideoLLM-online-8B-v1+ (Chen et al., 2024c)	63.1	_
VideoLLM-MoD (Wu et al., 2024)	63.4	_
VLog (Lin & Shou, 2025)	57.4	_
ProVideLLM-8B/11+ (Chatterjee et al., 2025)	67.3	_
VL-JEPA (FT)	77.2	86.9

#### 4.6 WORLDPREDICTION-WM

**Evaluation Setup.** We evaluate VL-JEPA on the "world modeling" task in the WORLDPREDICTION (Chen et al., 2025a) benchmark, where the model is provided with two images representing the initial and final world states and must identify, among four candidate video clips, the action that explains the observed transition. To adapt VL-JEPA, we duplicate and concatenate the initial and final state images to extract a *state embedding*, and encode each action candidate into *action embeddings*. The model then selects the candidate whose embedding is closest to the state embedding.

**Results.** Table 4 shows accuracy comparisons. VL-JEPA attains **61.8%** top-1 accuracy on WORLDPREDICTION-WM, establishing a new state of the art. Our VL-JEPA model not only substantially surpasses existing VLMs of comparable or larger scale but also exceeds the performance of frontier LLMs such as GPT-40, Claude-3.5-sonnet, and Gemini-2.0.

#### 4.7 ABLATION STUDY

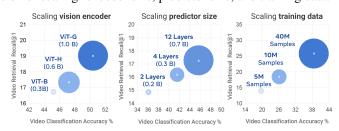
**Evaluation Setup.** We study different design choices for VL-JEPA. Each model variant is trained for 23k steps with a batch size of 256 on a dataset mixture of 3M samples. The mixture consists of 50% of video action recognition, 25% video captioning and 25% of image captioning, randomly

Table 4: WORLDPREDICTION-WM benchmark results. We compare the accuracy between large VLMs, socratic LLMs, and VL-JEPA. Our lightweight model achieves a new SoTA at 61.8%.

		Vision	n Lang	uage N	Aodels			Socratic LLMs (w/ Qwen2.5-VL-72B captions)							Ours			
	Inter	nVL2.	5		Qwen2	.5-VI		Llam	a-3.1	Llar	na-4	Q	wen2.	. 5	GPT-40	Claude-3.5	Gemini-2	VL-JEPA
2B	4B	26B	38B	3B	7в	32B	72B	8B	70B	109B	400B	3B	7в	72B	N/A	N/A	N/A	2B
20.0	29.8	30.2	50.3	21.6	45.5	49.0	<u>57.0</u>	48.7	<u>49.8</u>	52.7	<u>53.6</u>	44.0	49.1	48.5	52.0	<u>53.3</u>	<u>55.6</u>	61.8

Table 5: Ablation study results. Figure 5: Scalability Analysis. VL-JEPA benefits from 3 dimen-- sion of scaling: encoder size, predictor size, and training data.

Ablation	Setting	R@1	Acc.
Predictor	LLM (L8-16)	16.73	45.20
	LLM (L0-8)	15.80	35.86
	LLM (L4-12)	16.00	38.96
	LLM scratch	15.40	40.56
Attention	Bi-directional	21.13	49.82
	Causal	16.73	45.20
Loss	L2	16.73	45.20
	L1	15.20	36.22



sampled from our pretraining data. We evaluate on two task groups: average zero-shot text-to-video retrieval (Recall@1) on 4 datasets (validation split of PLM-Video-Auto, MSR-VTT, ActivityNet, DiDeMo), and average zero-shot close-vocabulary action recognition (top-1 accuracy) on 4 datasets (step and task recognition on COIN and CrossTask). We report the results in Tab. 5.

**Results.** Predictor initialization. Using mid-to-late LLM layers (L8--16) yields the best recognition accuracy (45.2), followed by L4–12 (38.9) and the early layers L0–8 (35.8). We observe a clear positive correlation between recognition performance and the depth of the LLM layers. Moreover, initializing from a pretrained LLM consistently improves recognition (LLM scratch consists of 8 layers). Predictor Attention Mask. Bi-directional attention outperforms causal attention on both retrieval (R@1 21.1 vs. 16.7) and recognition (49.8 vs. 45.2). This reinforces the previous observation that richer interactions between the query and visual tokens strengthen both tasks. Note that in this setting, the query token is appended to the end of the token sequence. **Training Loss.** L2 loss achieves the best overall balance, outperforming L1 for both retrieval (R@1 16.7 vs. 15.2) and recognition (45.2 vs. 36.2).

## 4.8 SCALABILITY ANALYSIS

Fig. 5 presents the results of scaling VL-JEPA. For model scaling, we use different sizes of V-JEPA 2 vision encoder, and vary the number of LLM Transformer layers for the predictor. For data scaling, we compare zero-shot mode performance with different training samples seen. Results show that both model scaling and data scaling yield consistent improvement in text-to-video retrieval and classification (i.e., video-to-text retrieval), measured on same group of datasets used in §4.7.

#### CONCLUSION

We have present VL-JEPA, a new vision-language model built upon the joint embedding predictive architecture. By shifting supervision from discrete token space to continuous semantic embedding space, VL-JEPA simplifies the learning target, avoids redundant modeling of surface linguistic variability, and enables non-autoregressive prediction. Through controlled experiments, we show that VL-JEPA outperforms generative VLMs trained with cross-entropy loss under matched training data budget, while achieving superior training efficiency and significantly lower inference latency. Beyond generation tasks, the embedding-based design further allows VL-JEPA to handle openvocabulary classification and cross-modal retrieval within a single unified architecture. Its ability to emit continuous semantic embeddings also makes it particularly well suited for real-time video applications, where selective decoding can improve both responsiveness and efficiency.

### REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Kumar Ashutosh, Santhosh Kumar Ramakrishnan, Triantafyllos Afouras, and Kristen Grauman. Video-mined task graphs for keystep recognition in instructional videos. *Advances in Neural Information Processing Systems*, 36:67833–67846, 2023.
- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15619–15629, 2023.
- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.
- Piyush Bagad, Makarand Tapaswi, and Cees GM Snoek. Test of time: Instilling video-language models with a sense of time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2503–2516, 2023.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Federico Baldassarre, Marc Szafraniec, Basile Terver, Vasil Khalidov, Francisco Massa, Yann Le-Cun, Patrick Labatut, Maximilian Seitzer, and Piotr Bojanowski. Back to the features: Dino as a foundation for video world models. *arXiv preprint arXiv:2507.19468*, 2025.
- Hritik Bansal, Yonatan Bitton, Idan Szpektor, Kai-Wei Chang, and Aditya Grover. Videocon: Robust video-language alignment via contrast captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13927–13937, 2024.
- Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv* preprint arXiv:2105.04906, 2021.
- Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mido Assran, and Nicolas Ballas. V-jepa: Latent video prediction for visual representation learning. 2023.
- Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Icml*, volume 2, pp. 4, 2021.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- Kevin Black, Manuel Y Galliker, and Sergey Levine. Real-time execution of action chunking flow policies. *arXiv preprint arXiv:2506.07339*, 2025.
- Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, et al. Perception encoder: The best visual embeddings are not at the output of the network. *arXiv preprint arXiv:2504.13181*, 2025.
- Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C Li, Adrien Bardes, Suzanne Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, et al. An introduction to vision-language modeling. *arXiv* preprint arXiv:2405.17247, 2024.
  - Qingqing Cao, Bhargavi Paranjape, and Hannaneh Hajishirzi. Pumer: Pruning and merging tokens for efficient vision language models. *arXiv preprint arXiv:2305.17530*, 2023.

- Dibyadip Chatterjee, Edoardo Remelli, Yale Song, Bugra Tekin, Abhay Mittal, Bharat Bhatnagar, Necati Cihan CamgÃkz, Shreyas Hampali, Eric Sauser, Shugao Ma, et al. Memory-efficient streaming videollms for real-time procedural video understanding. *arXiv preprint* arXiv:2504.13915, 2025.
  - Delong Chen, Zhao Wu, Fan Liu, Zaiquan Yang, Shaoqiu Zheng, Ying Tan, and Erjin Zhou. Protoclip: Prototypical contrastive language image pretraining. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
    - Delong Chen, Samuel Cahyawijaya, Jianfeng Liu, Baoyuan Wang, and Pascale Fung. Subobject-level image tokenization. *arXiv preprint arXiv:2402.14327*, 2024a.
    - Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan Wang. Visual instruction tuning with polite flamingo. In *Proceedings of the aaai conference on artificial intelligence*, volume 38, pp. 17745–17753, 2024b.
    - Delong Chen, Willy Chung, Yejin Bang, Ziwei Ji, and Pascale Fung. Worldprediction: A benchmark for high-level world modeling and long-horizon procedural planning. *arXiv* preprint *arXiv*:2506.04363, 2025a.
    - Delong Chen, Theo Moutakanni, Willy Chung, Yejin Bang, Ziwei Ji, Allen Bolourchi, and Pascale Fung. Planning with reasoning using vision language world model. *arXiv preprint arXiv:2509.02722*, 2025b.
    - Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18407–18418, 2024c.
    - Xi Chen, Xiao Wang, Soravit Changpinyo, Anthony J Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.
    - Jang Hyun Cho, Andrea Madotto, Effrosyni Mavroudi, Triantafyllos Afouras, Tushar Nagarajan, Muhammad Maaz, Yale Song, Tengyu Ma, Shuming Hu, Suyog Jain, et al. Perceptionlm: Openaccess data and models for detailed visual understanding. arXiv preprint arXiv:2504.13180, 2025.
    - Guangzhao Dai, Xiangbo Shu, Wenhao Wu, Rui Yan, and Jiachao Zhang. Gpt4ego: unleashing the potential of pre-trained models for zero-shot egocentric action recognition. *IEEE Transactions on Multimedia*, 2024.
    - Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. Advances in neural information processing systems, 36:49250–49267, 2023.
    - Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision*, 130(1):33–55, 2022.
    - Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 91–104, 2025.
    - Shangzhe Di, Zhelun Yu, Guanghao Zhang, Haoyuan Li, Tao Zhong, Hao Cheng, Bolin Li, Wanggui He, Fangxun Shu, and Hao Jiang. Streaming video question-answering with in-context video kycache retrieval. *arXiv preprint arXiv:2503.00540*, 2025.
    - Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. Sonar: sentence-level multimodal and language-agnostic representations. *arXiv* preprint arXiv:2308.11466, 2023.

- Zhengcong Fei, Mingyuan Fan, and Junshi Huang. A-jepa: Joint-embedding predictive architecture can listen. *arXiv preprint arXiv:2311.15830*, 2023.
  - Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter, Michal Klein, David Haldimann, Sai Aitharaju, Victor G Turrisi da Costa, Louis Béthune, Zhe Gan, et al. Multimodal autoregressive pre-training of large vision encoders. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 9641–9654, 2025.
  - Pascale Fung, Yoram Bachrach, Asli Celikyilmaz, Kamalika Chaudhuri, Delong Chen, Willy Chung, Emmanuel Dupoux, Hongyu Gong, Hervé Jégou, Alessandro Lazaric, et al. Embodied ai agents: Modeling the world. *arXiv preprint arXiv:2506.22355*, 2025.
  - Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36:27092–27112, 2023.
  - Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pp. 5842–5850, 2017.
  - Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18995–19012, 2022.
  - Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19383–19400, 2024.
  - Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pp. 4904–4916. PMLR, 2021.
  - Cijo Jose, Theo Moutakanni, Dahyun Kang, Federico Baldassarre, Timothee Darcet, Hu Xu, Daniel Li, Marc Szafraniec, Michael Ramamonjisoa, Maxime Oquab, Oriane Simeoni, Huy V. Vo, Patrick Labatut, and Piotr Bojanowski. DINOv2 Meets Text: A Unified Framework for Image- and Pixel-Level Vision-Language Alignment. In 2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 24905–24916, Los Alamitos, CA, USA, June 2025. IEEE Computer Society. doi: 10.1109/CVPR52734.2025.02319. URL https://doi.ieeecomputersociety.org/10.1109/CVPR52734.2025.02319.
  - Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
  - Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal inputs and outputs. In *International Conference on Machine Learning*, pp. 17283–17300. PMLR, 2023.
  - Yann LeCun. A path towards autonomous machine intelligence. Open Review, 62(1):1–62, 2022.
  - Hongyang Lei, Xiaolong Cheng, Qi Qin, Dan Wang, Huazhen Huang, Qingqing Gu, Yetao Wu, and Luo Ji. M3-jepa: Multimodal alignment via multi-gate moe based on the joint-embedding predictive architecture. In *Forty-second International Conference on Machine Learning*, 2025.
  - Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7331–7341, 2021.

- Xianhang Li, Haoqin Tu, Mude Hui, Zeyu Wang, Bingchen Zhao, Junfei Xiao, Sucheng Ren, Jieru Mei, Qing Liu, Huangjie Zheng, et al. What if we recaption billions of web images with llama-3? *arXiv preprint arXiv:2406.08478*, 2024.
  - Kevin Qinghong Lin and Mike Zheng Shou. Vlog: Video-language models by generative retrieval of narration vocabulary. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 3218–3228, 2025.
  - Kevin Qinghong Lin, Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Z Xu, Difei Gao, Rong-Cheng Tu, Wenzhe Zhao, Weijie Kong, et al. Egocentric video-language pretraining. *Advances in Neural Information Processing Systems*, 35:7575–7586, 2022a.
  - Xudong Lin, Fabio Petroni, Gedas Bertasius, Marcus Rohrbach, Shih-Fu Chang, and Lorenzo Torresani. Learning to recognize procedural activities with distant supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13853–13863, 2022b.
  - Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2024.
  - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
  - Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*, 2021.
  - Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvlm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025.
  - Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=8tYRqb05pVn.
  - Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2630–2640, 2019.
  - Fionn Murtagh and Pedro Contreras. Algorithms for hierarchical clustering: an overview. Wiley interdisciplinary reviews: data mining and knowledge discovery, 2(1):86–97, 2012.
  - Medhini Narasimhan, Licheng Yu, Sean Bell, Ning Zhang, and Trevor Darrell. Learning and verification of task structure in instructional videos. *arXiv preprint arXiv:2303.13519*, 2023.
  - Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
  - Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video understanding with large language models. *Advances in Neural Information Processing Systems*, 37:119336–119360, 2024.
  - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
  - Laura Sevilla-Lara, Shengxin Zha, Zhicheng Yan, Vedanuj Goswami, Matt Feiszli, and Lorenzo Torresani. Only time can tell: Discovering temporal data for temporal modeling. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 535–544, 2021.
  - Mustafa Shukor and Matthieu Cord. Skipping computations in multimodal llms. *arXiv preprint arXiv:2410.09454*, 2024.

- Mustafa Shukor, Corentin Dancette, and Matthieu Cord. ep-alm: Efficient perceptual augmentation of language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22056–22069, 2023.
- Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025.
- Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv* preprint arXiv:2508.10104, 2025.
- Wenxuan Song, Jiayi Chen, Pengxiang Ding, Han Zhao, Wei Zhao, Zhide Zhong, Zongyuan Ge, Jun Ma, and Haoang Li. Accelerating vision-language-action model integrated with action chunking via parallel decoding. *arXiv preprint arXiv:2503.02310*, 2025.
- Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1207–1216, 2019.
- Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- Théophane Vallaeys, Mustafa Shukor, Matthieu Cord, and Jakob Verbeek. Improved baselines for data-efficient perceptual augmentation of llms. In *European Conference on Computer Vision*, pp. 369–387. Springer, 2024.
- Pavan Kumar Anasosalu Vasu, Fartash Faghri, Chun-Liang Li, Cem Koc, Nate True, Albert Antony, Gokula Santhanam, James Gabriel, Peter Grasch, Oncel Tuzel, et al. Fastvlm: Efficient vision encoding for vision language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 19769–19780, 2025.
- Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.
- Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong Shi, et al. Internvideo2: Scaling foundation models for multimodal video understanding. In *European Conference on Computer Vision*, pp. 396–416. Springer, 2024.
- Shiwei Wu, Joya Chen, Kevin Qinghong Lin, Qimeng Wang, Yan Gao, Qianli Xu, Tong Xu, Yao Hu, Enhong Chen, and Mike Zheng Shou. Videollm-mod: Efficient video-language streaming with mixture-of-depths vision computation. *Advances in Neural Information Processing Systems*, 37:109922–109947, 2024.
- Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv* preprint arXiv:2109.14084, 2021.
- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data. *arXiv* preprint arXiv:2309.16671, 2023.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5288–5296, 2016.

- Shen Yan, Tao Zhu, Zirui Wang, Yuan Cao, Mi Zhang, Soham Ghosh, Yonghui Wu, and Jiahui Yu. Videococa: Video-text modeling with zero-shot transfer from contrastive captioners. *arXiv* preprint arXiv:2212.04979, 2022.
  - 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.
  - Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv* preprint *arXiv*:2205.01917, 2022.
  - Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 11975–11986, 2023.
  - Long Zhao, Nitesh B Gundavarapu, Liangzhe Yuan, Hao Zhou, Shen Yan, Jennifer J Sun, Luke Friedman, Rui Qian, Tobias Weyand, Yue Zhao, et al. Videoprism: A foundational visual encoder for video understanding. *arXiv preprint arXiv:2402.13217*, 2024.
  - Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. Learning video representations from large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6586–6597, 2023.
  - Yiwu Zhong, Licheng Yu, Yang Bai, Shangwen Li, Xueting Yan, and Yin Li. Learning procedure-aware video representation from instructional videos and their narrations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14825–14835, 2023.
  - Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. Dino-wm: World models on pre-trained visual features enable zero-shot planning. In *Forty-second International Conference on Machine Learning*, 2025.
  - Honglu Zhou, Roberto Martín-Martín, Mubbasir Kapadia, Silvio Savarese, and Juan Carlos Niebles. Procedure-aware pretraining for instructional video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10727–10738, 2023.
  - Luowei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
  - Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Nagrani, and Cordelia Schmid. Streaming dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18243–18252, 2024.
  - Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-task weakly supervised learning from instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3537–3545, 2019.

The appendix is organized as follows: in A we discuss the limitations of this work, in B we provide additional related work, and in C we present further retrieval results.

#### A LIMITATIONS

In this work, we demonstrated the advantages of VL-JEPA over standard VLMs, particularly in efficiency, streaming, and video—text tasks. Our goal at this stage, is not to propose a universal alternative to VLMs, as this would require broader evaluation on tasks such as reasoning, tool use, and agentic behaviors where current token generative VLMs excel. We focused on video understanding benchmarks like action recognition and captioning, though the architecture can naturally extend to tasks such as VQA with an adjusted training mixture. While our study considered only video and text modalities, the approach could be extended to others, such as audio, using dedicated encoders. We trained with an L2 loss to align predictions with text features, which assumes a unimodal target distribution and may limit performance on tasks with inherent multimodal target distribution. Finally, although our results show clear benefits from scaling parameters and dataset size, we did not fully explore this direction, leaving it for future work.

#### B ADDITIONAL RELATED WORKS

Efficient Vision Language Models. The growing size and training cost of VLMs has motivated efforts to improve efficiency. On the training side, strong performance can be achieved by updating only a subset of parameters, such as the vision–language connector (Tsimpoukelli et al., 2021; Alayrac et al., 2022; Vallaeys et al., 2024; Shukor et al., 2023; Koh et al., 2023; Merullo et al., 2023; Dai et al., 2023). At inference, efficiency is pursued through pruning parameters or visual tokens (Cao et al., 2023; Shukor & Cord, 2024; Vasu et al., 2025). For real-time use cases, recent work explores small VLMs (Yao et al., 2024; Marafioti et al., 2025) and heuristics to reduce query frequency in asynchronous inference (Shukor et al., 2025).

#### C ZERO-SHOT VIDEO RETRIEVAL RESULTS

We compare models on text-to-video retrieval under the Recall@1 (R@1) metric in Talbe 6. The "Data" column summarizes the total number of paired vision—language samples (image—text or video—text) exposed during all training stages. CLIP (Radford et al., 2021) uses 400M image-text pairs. CLIP4Clip (Luo et al., 2021) builds on CLIP with an additional 380k HowTo100M (Miech et al., 2019) finetuning clips. InternVideo (Wang et al., 2022)used CLIP pretraining and added 14.35M video-text pairs and 100M LAION image—text pairs. InternVideo2 (Wang et al., 2024) scales further with 100M video-text and 300M LAION image—text pairs. CoCa (Yu et al., 2022) uses 4.8B pairs from JFT-3B + ALIGN. VideoCoCa (Yan et al., 2022) extends CoCa pretraining with VideoCC3M. VideoPrism (Zhao et al., 2024) use pretrained CoCa (4.8B), WebLI (~1B), and 618M extra samples. Perception Encoder (PE) (Bolya et al., 2025) is trained on 2.3B MetaCLIP (Xu et al., 2023) data. SigLIP2 (Tschannen et al., 2025) is trained on 12B WebLI image-text pairs.

Our VL-JEPA is trained with only 0.05B data, which is orders of magnitude smaller than foundation models. Despite this, it achieves competitive results, notably outperforming VideoCoCa-B on YouCook2 while using  $\sim\!100\times$  less data. In addition, unlike other baselines we **do not use contrastive objective during training**. We emphasize that the goal of this paper is not to produce a state-of-the-art retrieval model, but rather to demonstrate that VL-JEPA achieves reasonable retrieval performance given its scale. Importantly, many VLMs cannot perform retrieval at all without architectural modifications.

## LLM USAGE

We used large language models (LLMs) solely as writing assistants for this paper. Specifically, they were employed to help rephrase sentences for clarity and readability. No content, ideas, or experimental results were generated by LLMs. The authors take full responsibility for the scientific contributions and all written content.

Table 6: **Text-to-video retrieval Recall@1**. "Data" column reports the total paired vision—language data (image—text or video—text) seen during training. Note that, unlike other baselines, VL-JEPA is not trained with contrastive objective and keep the text encoder frozen.

Model	Data	MSR-VTT	ActivityNet	DiDeMo	MSVD	YouCook2
VL-JEPA	0.06B	22.3	22.9	22.5	35.1	17.7
CLIP (Radford et al., 2021)	0.4B	30.6	-	-	36.2	-
CLIP4Clip (Luo et al., 2021)	0.4B	32.0	-	-	38.5	
InternVideo (Wang et al., 2022)	0.5B	40.7	30.7	31.5	43.4	-
InternVideo2-1B (Wang et al., 2024)	0.4B	51.9	60.4	57.0	58.1	-
InternVideo2-6B (Wang et al., 2024)	0.4B	55.9	63.2	57.9	59.3	-
PE-B (Bolya et al., 2025)	2.3B	47.6	39.0		50.4	-
PE-L (Bolya et al., 2025)	2.3B	50.3	46.4		57.2	-
PE-G (Bolya et al., 2025)	2.3B	51.2	54.7		59.7	-
CoCa-B (Yu et al., 2022)	4.8B	27.5	21.7	-	-	11.2
CoCa (Yu et al., 2022)	4.8B	30.0	28.5	-	-	16.8
VideoCoCa-B (Yan et al., 2022)	4.8B	31.2	29.6	-	-	16.5
VideoCoCa (Yan et al., 2022)	4.8B	34.3	34.5	-	-	20.3
VideoPrism-B (Zhao et al., 2024)	6.4B	37.0	49.6	-	_	-
VideoPrism-g (Zhao et al., 2024)	6.4B	39.7	52.7	-	_	-
SigLIP2-B (Tschannen et al., 2025)	12.0B	38.5	28.6	-	49.0	-
SigLIP2-G (Tschannen et al., 2025)	12.0B	43.1	38.3	-	54.3	