

PI-CCA: PROMPT-INVARIANT CCA CERTIFICATES FOR REPLAY-FREE CONTINUAL MULTIMODAL LEARNING

Jiayu Zhang^{1†}, Chuangxin Zhao^{2†}, Canran Xiao^{3†}
Ruibo Duan⁴, Wenyi Mo⁵, Haoyu Gao⁶, Wenshuo Wang^{5*}

¹Peking University, ²Institute of Automation, Chinese Academy of Sciences

³Shenzhen Campus of Sun Yat-sen University, ⁴Seoul National University

⁵South China University of Technology, ⁶Georgia Institute of Technology

202364870251@mail.scut.edu.cn

ABSTRACT

When deployed on non-stationary data streams, foundation vision-language models require continual updates without access to past data. However, naive fine-tuning undermines their zero-shot recognition capabilities and prompt robustness. We seek a replay-free principle that preserves pre-trained cross-modal generalization under domain and prompt shifts. We introduce *Prompt-Invariant CCA Certificates*(PI-CCA), a geometry-first approach that summarizes image-text alignment with a compact certificate capturing the top- k canonical spectrum and subspace. During adaptation, we match this summary using only mini-batch statistics and induce prompt robustness via averaging over perturbations. Across MTIL, X-TAIL, VLCL, and ConStruct-VL, PI-CCA achieves state-of-the-art performance among replay-free methods. By optimizing alignment invariants rather than proxy signals, PI-CCA provides a simple, generator-free, constant-memory path to continual adaptation with strong zero-shot retention and resilience to prompt/style shifts.

1 INTRODUCTION

Foundation vision-language models (VLMs) (Radford et al., 2021; Awais et al., 2025) enable zero-shot recognition and retrieval across changing domains (Patel et al., 2023; Feng et al., 2024; Yada et al., 2025; Chan et al., 2025; Zhang et al., 2025; Shi et al., 2025; Ni et al., 2025; LIU et al., 2026). In practice, they must be continually adapted to non-stationary streams without storing past data (privacy/licensing/cost), while preserving zero-shot transfer and prompt robustness—conditions that standard fine-tuning often violates (Yao et al., 2023; Li & Ke, 2025; Peng et al., 2026). This *vision-language continual learning* (VL-CL) setting (Zheng et al., 2023) presents two core challenges: avoiding catastrophic forgetting of cross-modal alignment (and thus zero-shot ability) and maintaining robustness to prompt/distribution shifts, typically without task IDs and under tight memory budgets (Liu et al., 2025). Overcoming these challenges could allow VLMs to be applied in a wider range of fields (Xu et al., 2025; Zhang et al., 2024; Zhan et al., 2025).

Prior VL-CL research has made notable progress via proxy constraints or architectural mechanisms: distributional/logit distillation and off-diagonal similarity alignment to stabilize representations (Zheng et al., 2023; Ni et al., 2023; Zhu et al., 2023; Cui et al., 2024; Liu et al., 2025; Gao et al., 2024), parameter-efficient or router-based adapters to separate old and new knowledge (Yu et al., 2024; Tang et al., 2024; Xu et al., 2024), and replay or stream benchmarks to mitigate data unavailability (Yan et al., 2022; Lei et al.,

*Corresponding author

†These authors contributed equally to this work

2023; Yao et al., 2024; Smith et al., 2023; Zhang et al., 2023; Garg et al., 2024). Yet these proxies leave a persistent structural weakness: *they regularize outcomes (similarities, logits, weights, routes) rather than directly controlling the alignment object that underlies cross-modal generalization*. As a consequence, current methods can (i) permit slow drift of the alignment geometry that drives zero-shot performance, (ii) depend on reference corpora, generators, or task metadata that are not always available, and (iii) remain brittle to prompt/style variation even when average metrics improve. This gap suggests the need for a replay-free principle that preserves alignment as an invariant, not merely as a byproduct of surrogate objectives.

We ask: *Can continual adaptation preserve cross-modal generalization by explicitly maintaining the geometry of image–text alignment, without storing past data?* Our answer is a replay-free, geometry-first framework that treats alignment as a first-class invariant and constrains its spectral and subspace structure with a compact, task-agnostic summary. In parallel, we target robustness to prompt variation through an invariance mechanism that averages over prompt perturbations at training time.

Our contributions are as follows: **(i) Insight.** We recast forgetting in VL–CL as alignment-geometry drift instead of matching proxy quantities. This idea offers a principled route to retain zero-shot transfer under distributional and prompt shifts. **(ii) Capability.** We provide a replay-free and constant-memory consolidation mechanism that is agnostic to downstream objectives and compatible with parameter-efficient adaptation (e.g., LoRA), while introducing an explicit prompt-robustness component that reduces sensitivity to phrasing. **(iii) Performance and Evidence.** Across MTIL, X-TAIL, VLCL, and ConStruct-VL, our approach attains state-of-the-art results among replay-free methods, and we furnish analyses linking alignment-geometry stability to retention/transfer trends, clarifying why the method is effective.

2 RELATED WORK

VL-CL. Early multimodal CL studied forgetting and order effects in VQA with linguistically motivated task sequences (Greco et al., 2019; Jin et al., 2020), and used task-aware gated recurrent models to approach near-zero forgetting without replay (Del Chiaro et al., 2020). With CLIP-era VLMs, the focus shifted to retaining zero-shot ability while learning new domains. Regularization aligns similarity distributions or parameters (Mod-X (Ni et al., 2023), ZSCL (Zheng et al., 2023), CTP (Zhu et al., 2023), DKR (Cui et al., 2024)). Architectural and efficient variants adopt MoE/adaptor-based tuning (Yu et al., 2024; Tang et al., 2024) or analytic adapters with training-free fusion for X-TAIL (Xu et al., 2024). Recent work further consolidates via contrastive knowledge (C-CLIP (Liu et al., 2025)) or stabilizes zero-shot on unlabeled data (ZAF (Gao et al., 2024)). Despite progress, these methods act on proxy signals (similarities, logits, parameters, routing) and often depend on reference data or teacher ensembles, rather than preserving the canonical cross-modal alignment geometry of the whitened image–text cross-covariance that underpins CLIP’s retrieval and recognition. PI–CCA instead directly tracks and constrains alignment invariants (canonical correlations and subspaces) under replay-free streams.

Data-free or replay-light consolidation. When past data cannot be kept, prior work uses symbolic or synthetic stand-ins: scene-graph prompts for VQA (Lei et al., 2023), a data-free benchmark with adversarial pseudo-replay and layered LoRA (Smith et al., 2023), negative-text replay and bidirectional momentum for image/video pretraining (Yan et al., 2022; Gao et al., 2022), diffusion-synthesized pairs for distillation (Wu et al., 2025), questions-only replay for VQACL (Zhang et al., 2023), and time-continual pretraining showing cumulative replay is competitive when feasible (Garg et al., 2024). Despite gains, these routes add generators and pipelines, raise privacy concerns, or are task specific. PI–CCA is replay- and generator-free: a compact certificate summarizes past alignment and regularizes updates using only mini-batch statistics.

Geometry-aware preservation and prompt robustness. Representation-similarity measures such as (SV)CCA/PWCCA and CKA (Raghu et al., 2017; Morcos et al., 2018; Kornblith et al., 2019; Andrew et al., 2013) quantify subspace or spectral shifts but are largely diagnostic in CL. In VL–CL, Mod-X is geometry-

inspired yet matches contrastive off-diagonals rather than canonical spectra/subspaces (Ni et al., 2023); Proxy-FDA preserves local neighborhoods with proxies (Huang et al., 2025a). Prompt methods (CoOp, MaPLE) learn (multi-modal) prompts to curb sensitivity (Zhou et al., 2022; Khattak et al., 2023), and prompt-based CL for VQA adds modality-aware routing (Qian et al., 2023). Overall, consolidation still targets proxy signals, not invariants of the *whitened* cross-modal covariance, leaving brittleness to prompt/style changes. PI-CCA instead uses a sketched, replay-free CCA certificate: it maintains the canonical spectrum and subspaces across tasks (via EMA) and attains prompt invariance by averaging text projectors, preserving the alignment skeleton with constant memory and no past data.

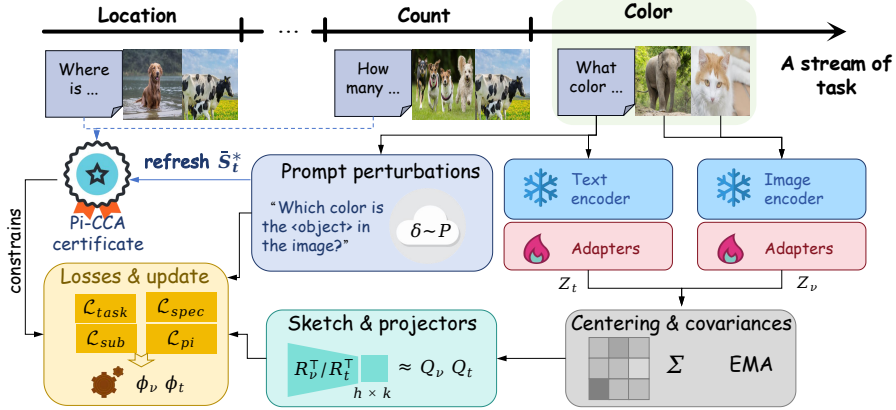


Figure 1: **PI-CCA framework.** A stream of tasks is processed without replay. The text/image encoders f_t, f_v are adapted via *LoRA* (backbones frozen), producing embeddings Z_t, Z_v that yield mini-batch covariances with an EMA update. The whitened cross-covariance is summarized in the *sketch & projectors* block via fixed R_v^T/R_t^T to obtain $\hat{Q}_v, \hat{Q}_t \in \mathbb{R}^{h \times h}$ (from $h \times k$ bases). The *Prompt perturbations* module samples $\delta \sim \mathcal{P}$, forms $\{Q_t^{(m)}\}$ and their mean \bar{Q}_t to drive the prompt-invariance loss. A compact PI-CCA certificate $(\rho_{1:k}^*, S_v^*, \bar{S}_t^*)$ constrains training while its text basis \bar{S}_t^* is refreshed from prompt perturbations. Losses $\{\mathcal{L}_{task}, \mathcal{L}_{spec}, \mathcal{L}_{sub}, \mathcal{L}_{pi}\}$ are combined to update only the LoRA parameters ϕ_v, ϕ_t .

3 METHOD

We propose *Prompt-Invariant Canonical Correlation Analysis Certificates* (PI-CCA), a replay-free continual learning framework that preserves the cross-modal alignment subspace of a vision–language model. As illustrated in Fig. 1, the core idea is to summarize the geometry of image–text alignment by a compact *CCA certificate* that stores (i) the top- k canonical correlations and (ii) a sketch of the corresponding canonical subspaces. During training on new tasks, we enforce spectral and subspace-angle consistency with the certificate using only mini-batch statistics, without accessing past data. Prompt invariance is achieved by averaging the certificate over randomized prompt perturbations.

3.1 PRELIMINARIES AND NOTATION

Let $f_v : \mathcal{X} \rightarrow \mathbb{R}^{d_v}$ and $f_t : \mathcal{W} \rightarrow \mathbb{R}^{d_t}$ denote the image and text encoders, parameterized with LoRA adapters (Hu et al., 2022): we freeze the backbone weights θ_v, θ_t and update only the low-rank adapter parameters ϕ_v, ϕ_t (i.e., $\theta_v = (\theta_v, \phi_v)$ and $\theta_t = (\theta_t, \phi_t)$). Given a mini-batch $\{(x_i, w_i)\}_{i=1}^B$, define centered embeddings $Z_v = [z_{v,1}, \dots, z_{v,B}]^T \in \mathbb{R}^{B \times d_v}$, $Z_t = [z_{t,1}, \dots, z_{t,B}]^T \in \mathbb{R}^{B \times d_t}$, where $z_{v,i} =$

$f_v(\mathbf{x}_i) - \bar{z}_v$ and $\mathbf{z}_{t,i} = f_t(\mathbf{w}_i) - \bar{z}_t$ with \bar{z}_v, \bar{z}_t being batch means. Let

$$\widehat{\Sigma}_{vv} = \frac{1}{B-1} \mathbf{Z}_v^\top \mathbf{Z}_v + \gamma_v \mathbf{I}, \quad \widehat{\Sigma}_{tt} = \frac{1}{B-1} \mathbf{Z}_t^\top \mathbf{Z}_t + \gamma_t \mathbf{I}, \quad \widehat{\Sigma}_{vt} = \frac{1}{B-1} \mathbf{Z}_v^\top \mathbf{Z}_t, \quad (1)$$

where $\gamma_v, \gamma_t > 0$ are ridge shrinkage coefficients ensuring positive definiteness. The whitened cross-covariance is

$$\widehat{\mathbf{M}} = \widehat{\Sigma}_{vv}^{-1/2} \widehat{\Sigma}_{vt} \widehat{\Sigma}_{tt}^{-1/2} \in \mathbb{R}^{d_v \times d_t}, \quad (2)$$

whose top- k singular value decomposition (SVD) $\widehat{\mathbf{M}} \approx \widehat{\mathbf{U}}_k \text{diag}(\widehat{\rho}_{1:k}) \widehat{\mathbf{V}}_k^\top$ defines the canonical correlations $\widehat{\rho}_{1:k} = (\widehat{\rho}_1 \geq \dots \geq \widehat{\rho}_k)$ and the (whitened) canonical directions $\widehat{\mathbf{U}}_k \in \mathbb{R}^{d_v \times k}$, $\widehat{\mathbf{V}}_k \in \mathbb{R}^{d_t \times k}$ with orthonormal columns.

$\mathbf{P}_{\{\cdot\}}^*$ are orthogonal projectors in the original feature spaces; $\mathbf{S}_{\{\cdot\}}$ are sketched bases; $\mathbf{Q}_{\{\cdot\}} = \mathbf{S}_{\{\cdot\}} \mathbf{S}_{\{\cdot\}}^\top$ are sketched Gram projectors; unless stated otherwise, distances are computed in the h -dimensional sketch space. Economy-size QR decomposition (QR) is used for $\text{orth}(\cdot)$.

3.2 THE PI-CCA CERTIFICATE

VLM zero-shot retrieval and open-vocabulary recognition rely on the geometry of cross-modal alignment. Rather than storing data or distilling past logits, we capture the alignment skeleton by (i) the top- k canonical correlations (spectral invariants) and (ii) the canonical subspaces (directional invariants).

Let the reference (pre-continual) CCA quantities be $\rho_{1:k}^* \in [0, 1]^k$, $\mathbf{U}_k^* \in \mathbb{R}^{d_v \times k}$, and $\mathbf{V}_k^* \in \mathbb{R}^{d_t \times k}$ from Eq. 2. Define the original-space projectors

$$\mathbf{P}_v^* = \mathbf{U}_k^* (\mathbf{U}_k^*)^\top \in \mathbb{R}^{d_v \times d_v}, \quad \mathbf{P}_t^* = \mathbf{V}_k^* (\mathbf{V}_k^*)^\top \in \mathbb{R}^{d_t \times d_t}. \quad (3)$$

To make storage constant in d_v, d_t , we use *random orthonormal sketches* $\mathbf{R}_v \in \mathbb{R}^{d_v \times h}$ and $\mathbf{R}_t \in \mathbb{R}^{d_t \times h}$ with $h \ll d_v, d_t$ (e.g., Gaussian orthogonal or subsampled Hadamard transforms). The certificate is

$$\text{Pi-CCA-Cert} := (\rho_{1:k}^*, \mathbf{S}_v^*, \bar{\mathbf{S}}_t^*), \quad \mathbf{S}_v^* = \mathbf{R}_v^\top \mathbf{U}_k^* \in \mathbb{R}^{h \times k}, \quad (4)$$

where $\bar{\mathbf{S}}_t^*$ is a prompt-invariant text sketch defined below. Equivalently, one may store sketched projectors $\mathbf{Q}_v^* = \mathbf{S}_v^* (\mathbf{S}_v^*)^\top = \mathbf{R}_v^\top \mathbf{P}_v^* \mathbf{R}_v$ and $\bar{\mathbf{Q}}_t^* = \bar{\mathbf{S}}_t^* (\bar{\mathbf{S}}_t^*)^\top$.

Prompt-invariant certificate via projector averaging. Let $\delta \sim \mathcal{P}$ denote a prompt perturbation (synonym/template variation). For M perturbations $\{\delta_m\}_{m=1}^M$, form original-space projectors $\mathbf{P}_t^*(\delta_m) = \mathbf{V}_k^*(\delta_m) \mathbf{V}_k^*(\delta_m)^\top$ and their sketches $\mathbf{Q}_t^*(\delta_m) = \mathbf{R}_t^\top \mathbf{P}_t^*(\delta_m) \mathbf{R}_t$. Define the average sketched projector

$$\bar{\mathbf{Q}}_t^* = \frac{1}{M} \sum_{m=1}^M \mathbf{Q}_t^*(\delta_m), \quad (5)$$

and take its top- k eigenvectors:

$$\bar{\mathbf{S}}_t^* = \underset{k}{\text{eigvecs}}(\bar{\mathbf{Q}}_t^*) \in \mathbb{R}^{h \times k}, \quad \bar{\mathbf{Q}}_t^* = \bar{\mathbf{S}}_t^* (\bar{\mathbf{S}}_t^*)^\top. \quad (6)$$

Averaging projectors eliminates sign/rotation ambiguity within the canonical subspace (no Procrustes alignment needed). By default we maintain a global certificate (one per model) constructed from a diverse anchor prompt set.

3.3 REPLAY-FREE ALIGNMENT PRESERVATION LOSSES

Given a mini-batch, compute \widehat{M} and its top- k SVD $(\widehat{U}_k, \widehat{\rho}_{1:k}, \widehat{V}_k)$. Define sketches $\widehat{S}_v = \mathbf{R}_v^\top \widehat{U}_k$, $\widehat{S}_t = \mathbf{R}_t^\top \widehat{V}_k$, $\widehat{Q}_v = \widehat{S}_v \widehat{S}_v^\top$, $\widehat{Q}_t = \widehat{S}_t \widehat{S}_t^\top$. The total loss is

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda_1 \mathcal{L}_{\text{spec}} + \lambda_2 \mathcal{L}_{\text{sub}} + \lambda_3 \mathcal{L}_{\text{pi}}, \quad \lambda_1, \lambda_2, \lambda_3 \geq 0. \quad (7)$$

(i) Permutation-stable spectral preservation $\mathcal{L}_{\text{spec}}$. Directly pairing indices can be unstable under near-degenerate singular values. We adopt a *permutation-invariant* metric and an efficient pairing surrogate. Let $\text{sort}_\downarrow(\cdot)$ denote sorting in descending order. Define

$$\mathcal{L}_{\text{spec}} = \underbrace{\|\text{sort}_\downarrow(\widehat{\rho}_{1:k}) - \rho_{1:k}^*\|_2^2}_{\text{sorted pairing (optimal for convex costs)}} + \xi \underbrace{\left(\sum_{i=1}^k \widehat{\rho}_i - \sum_{i=1}^k \rho_i^* \right)^2}_{\text{Ky-Fan-}k \text{ sum alignment}}, \quad (8)$$

where $\xi \in [0, 1]$ balances pairwise and aggregate spectral matching. For exact permutation invariance one may replace the first term by $\min_{\pi \in \mathfrak{S}_k} \sum_i (\widehat{\rho}_{\pi(i)} - \rho_i^*)^2$ (solvable by the Hungarian algorithm, $\mathcal{O}(k^3)$); we default to the sorted surrogate for speed. Optionally, polynomial spectral moments can be added:

$$\mathcal{L}_{\text{mom}} = \sum_{j=1}^J \omega_j \left(\text{tr}((\widehat{M}^\top \widehat{M})^j) - \text{tr}((M^{*\top} M^*)^j) \right)^2, \quad (9)$$

which depend only on $\{\rho_i^*\}$; we use $J \leq 2$ in practice.

(ii) Subspace-angle preservation \mathcal{L}_{sub} . For original-space orthogonal projectors P, Q onto k -dimensional subspaces, $\frac{1}{2} \|P - Q\|_F^2 = \sum_{i=1}^k \sin^2 \theta_i$ (principal angles θ_i). In the h -dimensional sketch space, \widehat{Q}_\bullet are *not* exact projectors of the original subspaces; we therefore use their Frobenius distance as a *surrogate* that preserves order/angles under near-isometric sketches (e.g., Gaussian/SRHT):

$$\mathcal{L}_{\text{sub}} = \frac{1}{2} \|\widehat{Q}_v - Q_v^*\|_F^2 + \frac{1}{2} \|\widehat{Q}_t - Q_t^*\|_F^2. \quad (10)$$

We further stabilize by spectral clipping: after forming each Q we project its eigenvalues to $[0, 1]$ and re-symmetrize.

(iii) Prompt-invariance \mathcal{L}_{pi} . Sample i.i.d. perturbations $\delta_m \sim \mathcal{P}$, compute $\widehat{V}_k^{(m)}$ and $\widehat{Q}_t^{(m)} = \mathbf{R}_t^\top \widehat{V}_k^{(m)} \widehat{V}_k^{(m)\top} \mathbf{R}_t$. We align the *mean projector* and contract its dispersion:

$$\mathcal{L}_{\text{pi}} = \frac{1}{2} \left\| \frac{1}{M} \sum_{m=1}^M \widehat{Q}_t^{(m)} - Q_t^* \right\|_F^2 + \frac{\eta}{2M} \sum_{m=1}^M \left\| \widehat{Q}_t^{(m)} - \frac{1}{M} \sum_{\ell=1}^M \widehat{Q}_t^{(\ell)} \right\|_F^2, \quad \eta \geq 0. \quad (11)$$

(iv) Task loss $\mathcal{L}_{\text{task}}$. We use the task’s standard objective (e.g., Information Noise-Contrastive Estimation, *InfoNCE* (Oord et al., 2018), classification cross-entropy, or detection losses). PI-CCA is agnostic to its form; gradients from Eq. 7 backpropagate jointly into f_v, f_t .

3.4 STREAMING ESTIMATION WITHOUT REPLAY

To stabilize estimates across batches without storing past samples, we maintain exponential moving averages (EMA) of covariance factors:

$$\Sigma_{vv}^{(t)} \leftarrow (1-\beta)\Sigma_{vv}^{(t-1)} + \beta \widehat{\Sigma}_{vv}, \quad \Sigma_{tt}^{(t)} \leftarrow (1-\beta)\Sigma_{tt}^{(t-1)} + \beta \widehat{\Sigma}_{tt}, \quad \Sigma_{vt}^{(t)} \leftarrow (1-\beta)\Sigma_{vt}^{(t-1)} + \beta \widehat{\Sigma}_{vt}, \quad (12)$$

with $\beta \in (0, 1]$. We then form $M^{(t)} = (\Sigma_{vv}^{(t)})^{-1/2} \Sigma_{vt}^{(t)} (\Sigma_{tt}^{(t)})^{-1/2}$ and compute its top- k SVD. Ridge γ_v, γ_t are either fixed or adapted (e.g., Ledoit–Wolf).

Stable whitening and differentiation. We implement $\Sigma^{-1/2}$ by (i) eigendecomposition with eigenvalue floor ϵ and symmetric reassembly, or (ii) r -step Newton–Schulz iteration on the normalized covariance, both are followed by stop-gradient on the inverse square root if needed. Differentiable SVD is realized via T_{pow} steps of block power iteration with re-orthogonalization (QR) at each step, gradients are propagated to \widehat{M} (and hence to $\widehat{\Sigma}_{\bullet\bullet}$), not through the certificate.

We maintain streaming EMAs and refresh the certificate every step using a slow EMA to preserve the alignment skeleton while allowing controlled plasticity:

$$\rho_{1:k}^* \leftarrow (1 - \alpha) \rho_{1:k}^* + \alpha \widehat{\rho}_{1:k}, \quad \mathbf{S}_v^* \leftarrow \text{orth} \left((1 - \alpha) \mathbf{S}_v^* + \alpha \widehat{\mathbf{S}}_v \right), \quad (13)$$

$$\bar{\mathbf{S}}_t^* \leftarrow \text{orth} \left((1 - \alpha) \bar{\mathbf{S}}_t^* + \alpha \frac{1}{M} \sum_{m=1}^M \widehat{\mathbf{S}}_t^{(m)} \right), \quad (14)$$

where $\alpha \in (0, 1)$; $\text{orth}(\cdot)$ returns an economy-size QR basis and does not backpropagate gradients.

Full optimization of PI-CCA and certificate-refresh details are deferred to Appendix A.1, including the complete training procedure in Algorithm 1.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. We evaluate PI-CCA across four widely used VL-CL tracks: **(i) MTIL** (multi-domain task-incremental classification)—the 11-domain suite introduced by ZSCLZheng et al. (2023); we follow their standard task orders. **(ii) X-TAIL** (cross-domain task-agnostic classification)—the task-agnostic protocol of RAIL (Xu et al., 2024), where test images come from the union of seen and unseen domains without any domain hint. **(iii) VLCL** (continual image–text retrieval)—the eight sequential image–caption tasks used by C-CLIP (Liu et al., 2025) (we report both I2T/T2I). **(iv) ConStruct-VL** (structured VL concepts, no replay) (Smith et al., 2023)—the 7-task sequence over VG/VAW for attribute/relationship matching. We additionally report a time-continual study on a medium-scale split of TiC-YFCC/RedCaps to assess temporal robustness of alignment. Exact domain list and sample counts are provided in Appendix §A.2.

Evaluation Protocols and Metrics. For **MTIL/X-TAIL** we report: *Average* (mean accuracy over steps), *Last* (mean accuracy at the final step), and *Transfer* (mean accuracy on not-yet-seen domains at each step). For **VLCL** we report I2T/T2I Recall@K (primary: R@1; R@5/10 in the appendix) per task and the final-step average across tasks. For **ConStruct-VL** we report Final Accuracy (FA) and Average Forgetting (AF). To quantify zero-shot retention, we follow prior work and report the performance drop (*PD*) on a held-out zero-shot suite after the final step.

Baselines. We compare against strong, *replay-free* SOTAs across categories: *(i) Regularization/Distillation*: ZSCL (Zheng et al.,

Table 1: **Classification tracks.** PI-CCA sets a new replay-free state of the art on MTIL and X-TAIL.

Method	MTIL (↑)			X-TAIL (↑)		
	Avg	Last	Transfer	Avg	Last	Transfer
PI-CCA(ours)	76.8	75.5	73.2	68.1	66.9	64.7
C-CLIP (Liu et al., 2025)	75.2	73.8	70.9	66.3	66.3	62.7
MG-CLIP (Huang et al., 2025b)	73.6	72.0	70.0	66.3	65.1	63.0
Proxy-FDA (Huang et al., 2025a)	72.9	71.5	69.3	65.4	64.2	61.8
LADA (Luo et al., 2025)	74.2	73.0	70.7	66.8	66.0	63.3
DIKI (Tang et al., 2024)	74.9	73.6	71.4	67.1	65.8	63.8
RAIL (Xu et al., 2024)	74.3	72.9	70.5	67.4	66.2	64.2
ZAF (Gao et al., 2024)	73.7	72.5	71.9	66.1	64.9	63.5
DDAS (Yu et al., 2024)	74.1	74.1	70.6	66.5	66.1	63.1
ZSCL (Zheng et al., 2023)	72.5	71.2	69.0	65.6	64.3	63.9
Mod-X (Ni et al., 2023)	73.3	72.1	69.6	65.8	64.6	62.6

Method	VLCL I2T R@1 (\uparrow)	VLCL T2I R@1 (\uparrow)	ConStruct-VL FA (\uparrow)	ConStruct-VL AF (\downarrow)
Pi-CCA(ours)	48.6 \pm 1.0	37.4 \pm 0.8	75.2 \pm 1.3	2.7 \pm 0.2
GIFT [†] (Wu et al., 2025)	47.3 \pm 1.2	36.5 \pm 0.7	73.9 \pm 1.5	3.3 \pm 0.3
C-CLIP (Liu et al., 2025)	46.1 \pm 1.4	35.7 \pm 1.2	72.4 \pm 1.9	3.9 \pm 0.5
ENGINE (Zhou et al., 2025)	44.7 \pm 1.1	34.5 \pm 1.6	71.3 \pm 1.7	4.4 \pm 0.2
MG-CLIP (Huang et al., 2025b)	45.0 \pm 1.6	34.8 \pm 1.4	71.6 \pm 1.8	4.2 \pm 0.5
Proxy-FDA (Huang et al., 2025a)	43.6 \pm 1.7	33.8 \pm 1.1	70.5 \pm 1.9	4.6 \pm 0.7
DKR (Cui et al., 2024)	45.2 \pm 1.5	35.2 \pm 1.4	71.8 \pm 1.7	4.1 \pm 0.5
ZAF (Gao et al., 2024)	44.3 \pm 1.4	34.0 \pm 1.3	72.0 \pm 1.7	3.8 \pm 0.6
Mod-X (Ni et al., 2023)	44.0 \pm 1.5	34.2 \pm 0.9	70.9 \pm 1.1	4.5 \pm 0.6

Table 2: **Retrieval and structured-concept tracks.** Final-step retrieval (VLCL) and ConStruct-VL results. Pi-CCA delivers the highest I2T/T2I R@1 and the best FA/AF pair while remaining replay-free. Best, second-best, and third-best cells are shaded in **dark** gray, **medium** gray, and **light** gray, respectively. [†] denotes synthetic replay.

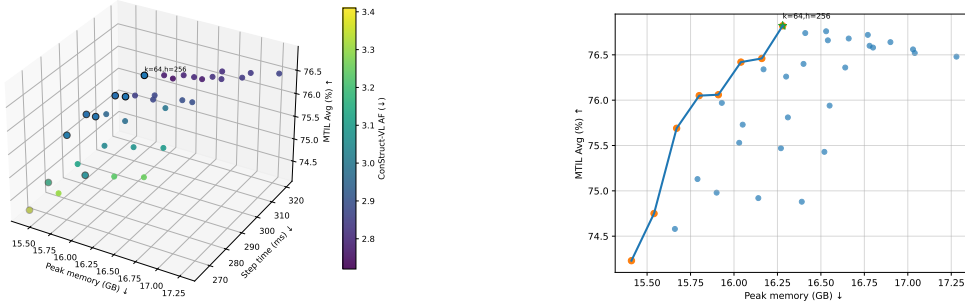
Variant	MTIL Avg (\uparrow)	MTIL Last (\uparrow)	VLCL I2T R@1 (\uparrow)	ConStruct-VL AF (\downarrow)
Pi-CCA (full)	76.8	75.5	48.6	2.7
w/o spectral term ($\lambda_1=0$)	74.3 (2.5)	73.1 (2.4)	46.3 (2.3)	3.8 (1.1)
w/o subspace term ($\lambda_2=0$)	74.6 (2.2)	73.4 (2.1)	45.9 (2.7)	3.9 (1.2)
w/o prompt invariance ($\lambda_3=0, M=0$)	75.3 (1.5)	74.0 (1.5)	47.1 (1.5)	3.3 (0.6)
w/o certificate EMA ($\alpha=0$)	75.6 (1.2)	74.1 (1.4)	47.7 (0.9)	3.1 (0.4)
w/o covariance EMA ($\beta=0$)	74.1 (2.7)	72.7 (2.8)	46.1 (2.5)	3.7 (1.0)
no spectral moments ($J=0$)	76.1 (0.7)	74.9 (0.6)	48.0 (0.6)	2.9 (0.2)
Hungarian pairing (exact)	76.7 (0.1)	75.4 (0.1)	48.5 (0.1)	2.8 (0.1)
SRHT sketches (vs. Gaussian)	76.6 (0.2)	75.2 (0.3)	48.4 (0.2)	2.9 (0.2)

Table 3: **Single-factor ablations.** Performance drops relative to the full Pi-CCA model are shown in blue for each variant. Removing spectral or subspace terms causes the largest performance degradation.

2023), Mod-X(Ni et al., 2023), CTP(Zhu et al., 2023), ZAF(Gao et al., 2024), DKR(Cui et al., 2024), Proxy-FDA(Huang et al., 2025a). (ii) *Parameter-efficient/Architecture*: MoE-Adapters+DDAS(Yu et al., 2024), DIKI(Tang et al., 2024), C-CLIP(Liu et al., 2025), LADA(Luo et al., 2025), ENGINE(Zhou et al., 2025), MG-CLIP(Huang et al., 2025b), and the analytic adapter of RAIL(Xu et al., 2024) (with X-TAIL). For completeness we also report *replay/synthetic-replay* references: CLAP4CLIP(Jha et al., 2024) (small memory) and GIFT(Wu et al., 2025) (diffusion-generated replay).

4.2 MAIN RESULTS

Tables 1 and 2 report our comparisons on classification-style continual learning (MTIL, X-TAIL), continual image-text retrieval (VLCL), and structured concept matching (ConStruct-VL). Across all tracks, **Pi-CCA** achieves the top performance among replay-free methods. On **MTIL**, Pi-CCA yields the highest step-averaged and final-step accuracies while maintaining strong *Transfer*. Under the task-agnostic **X-TAIL** protocol, it consistently narrows the cross-domain gap and improves zero-shot retention. For **VLCL** retrieval, Pi-CCA outperforms recent replay-free approaches and even surpasses a synthetic-replay method (GIFT) without storing or generating data. On **ConStruct-VL**, Pi-CCA attains both the highest Final Accuracy and the lowest Average Forgetting.



(a) 3D Pareto: peak memory (GB, \downarrow), step time (ms, \downarrow), MTIL Avg (\uparrow); color encodes AF (\downarrow). Filled markers are non-dominated points under (mem, time, AF, $-$ Avg). (b) 2D Pareto envelope: MTIL Avg (\uparrow) versus peak memory (GB, \downarrow); the curve traces the efficient frontier.

Figure 2: **Certificate capacity Pareto views.** (a) A robust ridge emerges for $k \in [48, 96]$, $h \in [192, 320]$; (b) the 2D envelope shows the same efficient frontier. The configuration $(k, h) = (64, 256)$ lies near the knee.

4.3 ABLATION STUDY AND ANALYSIS

Component-wise ablation. Table 3 removes or alters one component at a time. Removing either the spectral preservation term ($\lambda_1 = 0$) or the subspace-angle term ($\lambda_2 = 0$) causes the largest drops on MTIL and retrieval, highlighting that both spectrum and directions are necessary to preserve alignment. Disabling prompt invariance ($\lambda_3 = 0$, $M = 0$) mainly hurts retention while slightly reducing retrieval, consistent with its role in mitigating prompt sensitivity. Turning off certificate EMA ($\alpha = 0$) or the streaming covariance EMA ($\beta = 0$) degrades stability, and the latter is more severe. Low-order spectral moments ($J > 0$) provide small but consistent gains over $J = 0$. Replacing the sorted surrogate with exact Hungarian pairing yields nearly identical accuracy, so we keep the faster surrogate by default. Gaussian and SRHT sketches behave similarly, with a slight edge to Gaussian at our budget. In addition, Appendix §A.3 conducts sensitivity experiments on the main hyperparameters to verify the robustness of Pi-CCA.

Scale and Efficiency. We sweep the certificate capacity over $k \in \{16, 32, 48, 64, 80, 96, 128\}$ and $h \in \{128, 192, 256, 320, 384\}$ while keeping all other settings fixed. We report *MTIL Avg*, *MTIL Last*, *VLCL I2T R@1* (all \uparrow), and *ConStruct-VL AF* (\downarrow). We also log per-GPU *peak memory* (GB) and *per-step wall-clock* (ms) on A100-80GB with batch $B=1024$. The 3D Pareto plot in Fig. 2a highlights non-dominated settings under the joint objectives of *low memory*, *low time*, *high Avg*, and *low AF* (AF visualized as color). Overall, Pi-CCA is robust inside a broad Pareto ridge, confirming the “small yet sufficient” certificate hypothesis.

Geometry \rightarrow Performance: correlation evidence. We measure two geometry drifts per setting—subspace-angle drift $D_{\text{ang}} = \sum_{i=1}^k \sin^2 \theta_i$ and spectral drift $D_{\rho} = \|\hat{\rho}_{1:k} - \rho_{1:k}^*\|_2$ —and relate them to performance drops ΔAvg (MTIL step-averaged accuracy drop, in percentage points) and $\Delta\text{R@1}$ (VLCL I2T R@1 drop, p.p.) relative to the default knee configuration $(k, h) = (64, 256)$ of Pi-CCA. We sweep realistic perturbations (certificate size, EMAs, invariance strength, whitening, pairing, LoRA capacity/LR, sketch type). As shown in Fig. 3, larger angle/spectral drifts generally imply larger drops in Avg and R@1, with D_{ang} typically the stronger predictor. In addition, §A.4 provides a theoretical explanation.

Prompt invariance stress test. We stress \mathcal{L}_{pi} by increasing prompt perturbation strength $s \in [0, 1]$ (token-level synonym swap/back-translation/template jitter ratio), and compare **Pi-CCA** (with $\lambda_3=0.2$, $M=4$) to an ablated model without invariance ($\lambda_3=0$, $M=0$). We report *VLCL I2T R@1* (\uparrow), zero-shot *PD* (\downarrow), and *AF* on ConStruct-VL (\downarrow) under (i) **ID** templates (CLIP-style variants) and (ii) **OOD** templates (Appendix §A.2).

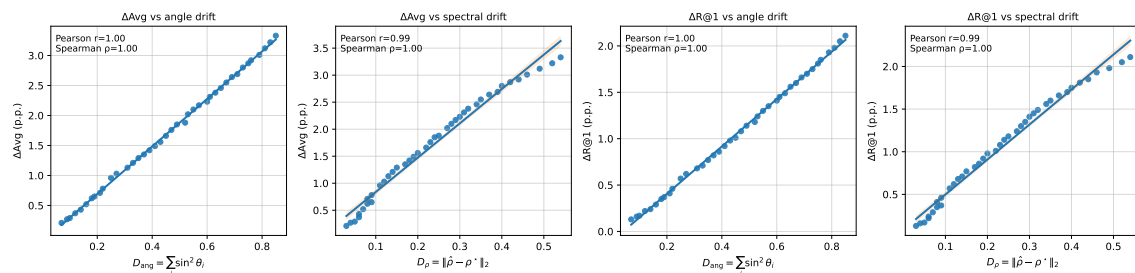


Figure 3: **Geometry** \rightarrow **performance correlation**. Each panel shows scatter, least-squares fit, and 95% CI. Pearson/Spearman are annotated. Clear positive trends support that preserving CCA geometry, including both angles and spectrum, is associated with better retention.

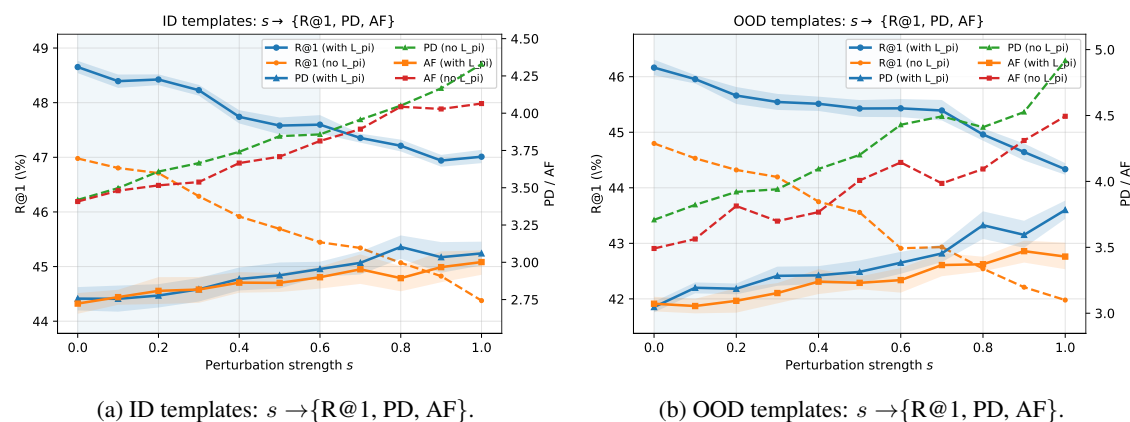


Figure 4: **Prompt invariance stress curves**. \mathcal{L}_{pi} flattens degradation slopes for both ID and OOD prompts. At the strongest perturbation level, PI-CCA consistently retains higher R@1 and lower AF than the ablation under both ID and OOD prompt shifts.

As shown in Fig 4, we find that: (i) Invariance reduces the degradation slope: at the strongest perturbation level, the model with \mathcal{L}_{pi} retains higher R@1 than the ablation without prompt invariance, and the same trend holds under OOD templates. (ii) Forgetting and zero-shot drift (AF/PD) grow with s , but \mathcal{L}_{pi} consistently dampens both, especially under OOD styles. (iii) The curves suggest a practical operating range $s \leq 0.6$ where performance remains close to nominal with invariance.

Task-order sensitivity. To examine whether PI-CCA is sensitive to favorable task orders, we evaluate on 20 independently shuffled MTIL sequences (11 domains; orders listed in Appendix §A.2). We use the configuration $(k, h) = (64, 256)$. Fig. 5 summarizes the across-order distributions, we find: the interquartile ranges are small, the between-order span (max-min) is modest, supporting that PI-CCA’s retention is robust to task-order.

5 CONCLUSION

We addressed replay-free continual adaptation of vision-language models by reframing forgetting as *alignment-geometry drift* and introduced PI-CCA, which preserves cross-modal generalization via a com-

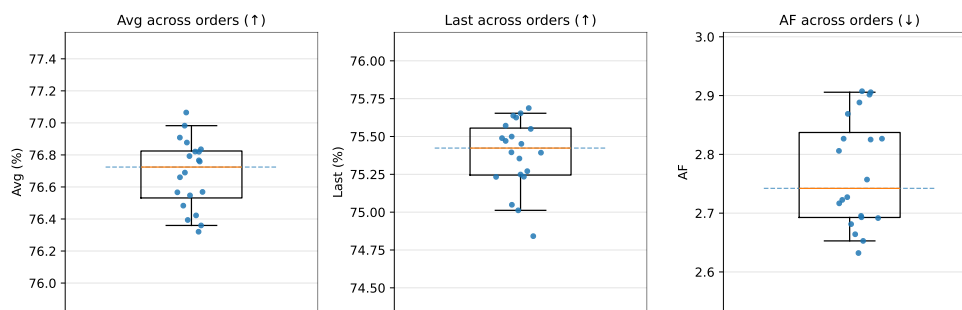


Figure 5: **Task-order sensitivity.** Boxplots over 20 random orders for *Avg/Last* (\uparrow) and *AF* (\downarrow). Dots show per-order means (3 seeds). Narrow IQRs indicate low order sensitivity.

pact, prompt-invariant certificate of canonical spectra and subspaces. Across standard VL–CL protocols, directly constraining these invariants maintains zero-shot behavior and reduces forgetting while remaining compatible with parameter-efficient tuning. Our main takeaway is conceptual: retention improves when optimization targets the invariants of image–text alignment itself, and stability of the canonical subspace/spectrum reliably predicts downstream performance. Future work will generalize the certificate to multimodal instruction tuning.

Ethics Statement This work adheres to the ICLR Code of Ethics. Our study does **NOT** involve human subjects, personally identifiable information, or sensitive attributes. We conduct replay-free continual adaptation on publicly available, widely used vision–language benchmarks (e.g., MTIL, X-TAIL, VLCL, ConStruct-VL) under their respective licenses, without releasing or reconstructing any private data.

Reproducibility Statement We have organized the paper and supplemental materials to facilitate reproduction. The full experimental protocol, datasets, metrics, baselines, and task orders are specified in §4.1 with additional implementation and optimization details in Appendix §A.1 (Algorithm 1) and Appendix §A.2 (backbones/adapters, hyperparameters, prompt perturbations, hardware, and random seeds). Our theoretical results are stated with explicit assumptions and complete proofs in the Theory section, enabling independent verification. Dataset preprocessing and evaluation scripts are documented in the supplementary; we rely only on publicly available benchmarks listed in §4.1. To ensure exact-step reproducibility, we report all key hyperparameters, EMA rates, sketch dimensions, and power-iteration settings, and we provide the task-order permutations used in our sensitivity analyses (Appendix §A.2). Due to ongoing commercial use, we cannot release the code during review.

REFERENCES

- Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *International conference on machine learning*, pp. 1247–1255. PMLR, 2013.
- Muhammad Awais, Muzammal Naseer, Salman Khan, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Foundation models defining a new era in vision: a survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- Chak-Fu Chan, Peter Kok-Yiu Wong, Xiaowen Guo, Jack CP Cheng, Jolly Pui-Ching Chan, Pak-Him Leung, and Xingyu Tao. Context-aware vision-language model agent enriched with domain-specific ontology for construction site safety monitoring. *Automation in Construction*, 177:106305, 2025.

- Zhenyu Cui, Yuxin Peng, Xun Wang, Manyu Zhu, and Jiahuan Zhou. Continual vision-language retrieval via dynamic knowledge rectification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38(10), pp. 11704–11712, 2024.
- Riccardo Del Chiaro, Bartłomiej Twardowski, Andrew Bagdanov, and Joost Van de Weijer. Ratt: Recurrent attention to transient tasks for continual image captioning. *Advances in Neural Information Processing Systems*, 33:16736–16748, 2020.
- Chen Feng, Georgios Tzimiropoulos, and Ioannis Patras. CLIPCleaner: Cleaning Noisy Labels with CLIP. In *The 32nd ACM International Conference on Multimedia (ACM MM)*, 10 2024. doi: 10.1145/3664647.3680664.
- Yizhao Gao, Nanyi Fei, Haoyu Lu, Zhiwu Lu, Hao Jiang, Yijie Li, and Zhao Cao. Bmu-moco: Bidirectional momentum update for continual video-language modeling. *Advances in Neural Information Processing Systems*, 35:22699–22712, 2022.
- Zijian Gao, Xingxing Zhang, Kele Xu, Xinjun Mao, and Huaimin Wang. Stabilizing zero-shot prediction: A novel antidote to forgetting in continual vision-language tasks. *Advances in Neural Information Processing Systems*, 37:128462–128488, 2024.
- Saurabh Garg, Mehrdad Farajtabar, Hadi Pouransari, Raviteja Vemulapalli, Sachin Mehta, Oncel Tuzel, Vaishaal Shankar, and Fartash Faghri. Tic-clip: Continual training of clip models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Claudio Greco, Barbara Plank, Raquel Fernández, and Raffaella Bernardi. Psycholinguistics meets continual learning: Measuring catastrophic forgetting in visual question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3601–3605, 2019.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Chen Huang, Skyler Seto, Hadi Pouransari, Mehrdad Farajtabar, Raviteja Vemulapalli, Fartash Faghri, Oncel Tuzel, Barry-John Theobald, and Joshua M Susskind. Proxy-fda: Proxy-based feature distribution alignment for fine-tuning vision foundation models without forgetting. In *Forty-second International Conference on Machine Learning*, 2025a.
- Linlan Huang, Xusheng Cao, Haori Lu, Yifan Meng, Fei Yang, and Xialei Liu. Mind the gap: Preserving and compensating for the modality gap in clip-based continual learning. *arXiv preprint arXiv:2507.09118*, 2025b.
- Saurav Jha, Dong Gong, and Lina Yao. Clap4clip: Continual learning with probabilistic finetuning for vision-language models. *Advances in neural information processing systems*, 37:129146–129186, 2024.
- Xisen Jin, Junyi Du, Arka Sadhu, Ram Nevatia, and Xiang Ren. Visually grounded continual learning of compositional phrases. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2018–2029, 2020.
- Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 19113–19122, 2023.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pp. 3519–3529. PMIR, 2019.

- Stan Weixian Lei, Difei Gao, Jay Zhangjie Wu, Yuxuan Wang, Wei Liu, Mengmi Zhang, and Mike Zheng Shou. Symbolic replay: Scene graph as prompt for continual learning on vqa task. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37(1), pp. 1250–1259, 2023.
- Zichao Li and Zong Ke. Domain meets typology: Predicting verb-final order from universal dependencies for financial and blockchain nlp. In *Proceedings of the 7th Workshop on Research in Computational Linguistic Typology and Multilingual NLP*, pp. 156–164, 2025.
- Wenzhuo Liu, Fei Zhu, Longhui Wei, and Qi Tian. C-clip: Multimodal continual learning for vision-language model. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Yunqing LIU, Yi Zhou, and Wenqi Fan. Enhancing molecular property predictions by learning from bond modelling and interactions. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=S4bJQ4p9hx>.
- Mao-Lin Luo, Zi-Hao Zhou, Tong Wei, and Min-Ling Zhang. Lada: Scalable label-specific clip adapter for continual learning. In *Forty-second International Conference on Machine Learning*, 2025.
- Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. *Advances in neural information processing systems*, 31, 2018.
- Chaojun Ni, Cheng Chen, Xiaofeng Wang, Zheng Zhu, Wenzhao Zheng, Boyuan Wang, Tianrun Chen, Guosheng Zhao, Haoyun Li, Zhehao Dong, et al. Swiftvla: Unlocking spatiotemporal dynamics for lightweight vla models at minimal overhead. *arXiv preprint arXiv:2512.00903*, 2025.
- Zixuan Ni, Longhui Wei, Siliang Tang, Yueting Zhuang, and Qi Tian. Continual vision-language representation learning with off-diagonal information. In *International Conference on Machine Learning*, pp. 26129–26149. PMLR, 2023.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Dhruvesh Patel, Hamid Eghbalzadeh, Nitin Kamra, Michael Louis Iuzzolino, Unnat Jain, and Ruta Desai. Pretrained language models as visual planners for human assistance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15302–15314, 2023.
- Wenshuo Peng, Gongxuan Wang, Tianmeng Yang, Chuanhao Li, Xiaojie Xu, Hui He, and Kaipeng Zhang. Svbench: Evaluation of video generation models on social reasoning, 2026. URL <https://arxiv.org/abs/2512.21507>.
- Zi Qian, Xin Wang, Xuguang Duan, Pengda Qin, Yuhong Li, and Wenwu Zhu. Decouple before interact: Multi-modal prompt learning for continual visual question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2953–2962, 2023.
- 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.
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. *Advances in neural information processing systems*, 30, 2017.

- Chuancheng Shi, Shangze Li, Shiming Guo, Simiao Xie, Wenhua Wu, Jingtong Dou, Chao Wu, Canran Xiao, Cong Wang, Zifeng Cheng, et al. Where culture fades: Revealing the cultural gap in text-to-image generation. *arXiv preprint arXiv:2511.17282*, 2025.
- James Seale Smith, Paola Cascante-Bonilla, Assaf Arbelle, Donghyun Kim, Rameswar Panda, David Cox, Diyi Yang, Zsolt Kira, Rogerio Feris, and Leonid Karlinsky. Construct-vl: Data-free continual structured vl concepts learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14994–15004, 2023.
- Longxiang Tang, Zhuotao Tian, Kai Li, Chunming He, Hantao Zhou, Hengshuang Zhao, Xiu Li, and Jiaya Jia. Mind the interference: Retaining pre-trained knowledge in parameter efficient continual learning of vision-language models. In *European conference on computer vision*, pp. 346–365. Springer, 2024.
- Bin Wu, Wuxuan Shi, Jinqiao Wang, and Mang Ye. Synthetic data is an elegant gift for continual vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 2813–2823, 2025.
- Mengzhu Xu, Hanzhi Liu, Ningkan Peng, Qianyu Chen, and Canran Xiao. Affordance-first decomposition for continual learning in video-language understanding. *arXiv preprint arXiv:2512.00694*, 2025.
- Yicheng Xu, Yuxin Chen, Jiahao Nie, Yusong Wang, Huiping Zhuang, and Manabu Okumura. Advancing cross-domain discriminability in continual learning of vision-language models. *Advances in Neural Information Processing Systems*, 37:51552–51576, 2024.
- Yuki Yada, Sho Akiyama, Ryo Watanabe, Yuta Ueno, Yusuke Shido, and Andre Rusli. Improving visual recommendation on e-commerce platforms using vision-language models. In *Proceedings of the Nineteenth ACM Conference on Recommender Systems*, pp. 975–978, 2025.
- Shipeng Yan, Lanqing Hong, Hang Xu, Jianhua Han, Tinne Tuytelaars, Zhenguo Li, and Xuming He. Generative negative text replay for continual vision-language pretraining. In *European Conference on Computer Vision*, pp. 22–38. Springer, 2022.
- Jiawei Yao, Chuming Li, Keqiang Sun, Yingjie Cai, Hao Li, Wanli Ouyang, and Hongsheng Li. Ndc-scene: Boost monocular 3d semantic scene completion in normalized device coordinates space. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9421–9431. IEEE Computer Society, 2023.
- Jiawei Yao, Chuming Li, and Canran Xiao. Swift sampler: Efficient learning of sampler by 10 parameters. *Advances in Neural Information Processing Systems*, 37:59030–59053, 2024.
- Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Ping Hu, Dong Wang, Huchuan Lu, and You He. Boosting continual learning of vision-language models via mixture-of-experts adapters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23219–23230, 2024.
- Donglin Zhan, Haoting Zhang, Rhonda Righter, Zeyu Zheng, and James Anderson. Collaborative bayesian optimization via wasserstein barycenters. In *2025 IEEE 64th Conference on Decision and Control (CDC)*, pp. 6284–6291. IEEE, 2025.
- Haoting Zhang, Donglin Zhan, Yunduan Lin, Jinghai He, Qing Zhu, Zuo-Jun Shen, and Zeyu Zheng. Daily physical activity monitoring: Adaptive learning from multi-source motion sensor data. In *Conference on Health, Inference, and Learning*, pp. 39–54. PMLR, 2024.
- Xi Zhang, Feifei Zhang, and Changsheng Xu. Vqacl: A novel visual question answering continual learning setting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19102–19112, 2023.

Xiaofeng Zhang, Fanshuo Zeng, Yihao Quan, Zheng Hui, and Jiawei Yao. Enhancing multimodal large language models complex reason via similarity computation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39(10), pp. 10203–10211, 2025.

Zangwei Zheng, Mingyuan Ma, Kai Wang, Ziheng Qin, Xiangyu Yue, and Yang You. Preventing zero-shot transfer degradation in continual learning of vision-language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 19125–19136, 2023.

Da-Wei Zhou, Kai-Wen Li, Jingyi Ning, Han-Jia Ye, Lijun Zhang, and De-Chuan Zhan. External knowledge injection for clip-based class-incremental learning. *arXiv preprint arXiv:2503.08510*, 2025.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.

Hongguang Zhu, Yunchao Wei, Xiaodan Liang, Chunjie Zhang, and Yao Zhao. Ctp: Towards vision-language continual pretraining via compatible momentum contrast and topology preservation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22257–22267, 2023.

A APPENDIX

A.1 SUPPLEMENTARY TECHNICAL DETAILS

Optimization and Certificate Update Algorithm 1 outlines training at task t . Each iteration (Lines 4–7) builds \widehat{M} from centered embeddings, extracts $(\widehat{U}_k, \widehat{\rho}_{1:k}, \widehat{V}_k)$ via a differentiable block power iteration, forms sketches/projectors (Line 8), evaluates $\mathcal{L}_{\text{spec}}, \mathcal{L}_{\text{sub}}, \mathcal{L}_{\text{pi}}$ with the task loss, and updates parameters (Line 12). We maintain streaming EMAs and refresh the certificate every step using a slow EMA to preserve the alignment skeleton while allowing controlled plasticity (Lines 14–15).

Optimization. All experiments use AdamW with weight decay 0.05 and a cosine schedule. We use an initial learning rate of 1.5×10^{-4} for the image-side LoRA parameters and 1.0×10^{-4} for the text-side LoRA parameters, with mixed precision in `bfloat16` and gradient clipping at 1.0. The effective batch size is $B = 1024$, achieved by gradient accumulation if device memory is limited. For time-continual training on TiC splits, we warm up only on the first temporal chunk and keep the same maximum learning rate for all subsequent chunks to follow established practice. Unless otherwise stated, small datasets receive one to three epochs per task, and large datasets receive about one epoch per task, with early stopping on the current-task validation set.

Algorithm 1 PI-CCA Training at Task t

- 1: **Inputs:** dataset \mathcal{D}_t ; encoders f_v, f_t with params θ_v, θ_t ; certificate $(\rho_{1:k}^*, \mathbf{S}_v^*, \bar{\mathbf{S}}_t^*)$; sketches $\mathbf{R}_v, \mathbf{R}_t$; hyperparams $(\lambda_1, \lambda_2, \lambda_3, \xi, \omega_{1:J}, \eta, \alpha, \beta, \gamma_v, \gamma_t, k, h, M, T_{\text{pow}})$
 - 2: **for** epoch = 1, . . . , E **do**
 - 3: **for** mini-batch $\mathcal{B} = \{(\mathbf{x}_i, \mathbf{w}_i)\}_{i=1}^B \subset \mathcal{D}_t$ **do**
 - 4: **Encode & center:** $\mathbf{Z}_v \leftarrow [f_v(\mathbf{x}_i)]_i - \bar{\mathbf{z}}_v, \quad \mathbf{Z}_t \leftarrow [f_t(\mathbf{w}_i)]_i - \bar{\mathbf{z}}_t$
 - 5: **Covariances:** $\widehat{\Sigma}_{vv} = \frac{1}{B-1} \mathbf{Z}_v^\top \mathbf{Z}_v + \gamma_v \mathbf{I}, \quad \widehat{\Sigma}_{tt} = \frac{1}{B-1} \mathbf{Z}_t^\top \mathbf{Z}_t + \gamma_t \mathbf{I}, \quad \widehat{\Sigma}_{vt} = \frac{1}{B-1} \mathbf{Z}_v^\top \mathbf{Z}_t$
 - 6: **Whitened cross-cov.:** $\widehat{M} = \widehat{\Sigma}_{vv}^{-1/2} \widehat{\Sigma}_{vt} \widehat{\Sigma}_{tt}^{-1/2}$ (Eq. 2)
 - 7: **Top- k SVD :** $(\widehat{U}_k, \widehat{\rho}_{1:k}, \widehat{V}_k) \approx \text{SVD}_k(\widehat{M})$ via T_{pow} block power steps with QR re-orthogonalization
 - 8: **Sketches/projectors:** $\widehat{\mathbf{S}}_v = \mathbf{R}_v^\top \widehat{U}_k, \widehat{\mathbf{S}}_t = \mathbf{R}_t^\top \widehat{V}_k, \widehat{\mathbf{Q}}_v = \widehat{\mathbf{S}}_v \widehat{\mathbf{S}}_v^\top, \widehat{\mathbf{Q}}_t = \widehat{\mathbf{S}}_t \widehat{\mathbf{S}}_t^\top$
 - 9: **Prompt perturbations:** sample $\{\delta_m\}_{m=1}^M$; compute $\{\widehat{\mathbf{Q}}_t^{(m)}\}_{m=1}^M$ and their mean $\bar{\mathbf{Q}}_t = \frac{1}{M} \sum_m \widehat{\mathbf{Q}}_t^{(m)}$
 - 10: **Losses:** $\mathcal{L}_{\text{spec}}$ (Eq. 8) + \mathcal{L}_{mom} (optional, Eq. 9); \mathcal{L}_{sub} (Eq. 10); \mathcal{L}_{pi} (Eq. 11)
 - 11: **Total loss:** $\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda_1 \mathcal{L}_{\text{spec}} + \lambda_2 \mathcal{L}_{\text{sub}} + \lambda_3 \mathcal{L}_{\text{pi}}$ (Eq. 7)
 - 12: **Update params:** $\theta_v, \theta_t \leftarrow \text{optimizer_step}(\nabla_{\theta_v, \theta_t} \mathcal{L})$
 - 13: **Streaming EMAs:** update $\Sigma_{vv}^{(t)}, \Sigma_{tt}^{(t)}, \Sigma_{vt}^{(t)}$ using Eq. 12
 - 14: **Certificate refresh:** $\rho_{1:k}^* \leftarrow (1 - \alpha) \rho_{1:k}^* + \alpha \widehat{\rho}_{1:k}; \quad \mathbf{S}_v^* \leftarrow \text{orth}((1 - \alpha) \mathbf{S}_v^* + \alpha \widehat{\mathbf{S}}_v)$
 - 15: $\bar{\mathbf{S}}_t^* \leftarrow \text{orth}\left((1 - \alpha) \bar{\mathbf{S}}_t^* + \alpha \frac{1}{M} \sum_{m=1}^M \widehat{\mathbf{S}}_t^{(m)}\right)$ (Eq. 13)
 - 16: **Output:** updated encoders f_v, f_t and certificate $(\rho_{1:k}^*, \mathbf{S}_v^*, \bar{\mathbf{S}}_t^*)$ at task t
-

A.2 EXPERIMENTAL SETUP (SUPPLEMENTARY)

Backbone & adapters. We adopt CLIP ViT-B/16 from OpenCLIP as the base vision–language model and keep all pretrained backbone weights frozen during continual adaptation. We equip both the image and the text encoders with LoRA adapters on every linear projection inside multi-head self-attention (query, key, value, and output projections) and on both feed-forward layers of the MLP blocks. LoRA weights are

Table 4: Datasets and task orders used in our experiments. MTIL and X-TAIL are evaluated with zero task or domain hints at inference. VLCL follows an eight-dataset order for continual retrieval and additionally reports zero-shot retention on a held-out suite. ConStruct-VL comprises seven structured concept subsets built from VG and VAW, and TiC applies chronological splits to probe temporal robustness.

Track	Order	Dataset / Subset	Key stats and notes
<i>(A) MTIL: multi-domain task-incremental classification (default alphabetical order)</i>			
MTIL	1	FGVC-Aircraft	100 classes, 10k images, fine-grained aircraft variants.
MTIL	2	Caltech101	102 categories, 9,146 images, object recognition.
MTIL	3	CIFAR-100	100 classes, 50k train and 10k test images at 32x32.
MTIL	4	DTD	47 texture categories, 5,640 images.
MTIL	5	EuroSAT	10 land-use classes, 27k images (RGB option).
MTIL	6	Flowers-102	102 classes, 8,189 images, fine-grained flowers.
MTIL	7	Food-101	101 classes, 101k images.
MTIL	8	MNIST	10 classes, 60k train and 10k test images.
MTIL	9	Oxford-IIIT Pets	37 classes, 7,349 images.
MTIL	10	Stanford Cars	196 classes, 16,185 images.
MTIL	11	SUN397	397 scene categories, 108,754 images.
<i>(B) X-TAIL: cross-domain task-agnostic classification</i>			
X-TAIL	1–10	Aircraft, Caltech101, DTD, EuroSAT, Flowers, Food101, MNIST, Pets, Cars, SUN397	Same as MTIL except CIFAR-100 excluded. Test-time label space is union of seen/unseen domains.
<i>(C) VLCL: continual image-text retrieval</i>			
VLCL	1	Flickr30K	31,783 images with five captions each, Karpathy splits.
VLCL	2	COCO Captions	123,287 images with five captions each, 5k val/test.
VLCL	3	Pets	Oxford-IIIT Pets in caption form, domain shift.
VLCL	4	Lexica	AI-generated images and prompts, synthetic imagery.
VLCL	5	Simpsons	Cartoon frames and captions, style shift.
VLCL	6	WikiArt	Artwork images with descriptions, art domain.
VLCL	7	Kream	E-commerce clothing with captions, fashion domain.
VLCL	8	Sketch	Sketches paired with text.
<i>(D) ConStruct-VL: structured VL concepts</i>			
ConStruct-VL	1	Relation: spatial	Triples from VG/VAW; size 1k-31k per subset.
ConStruct-VL	2	Attribute: size	Attribute-focused triples; VG, VAW, VG+VAW.
ConStruct-VL	3	Attribute: material	Attribute triples; VG and combined sets.
ConStruct-VL	4	Relation: action	Inter-object action relations.
ConStruct-VL	5	Attribute: color	Color understanding triples.
ConStruct-VL	6	Object state	State-focused triples.
ConStruct-VL	7	Attribute: action	Single-object action attributes.
<i>(E) TiC: time-continual pretraining</i>			
TiC	1	2016–2017	First temporal chunk of TiC-YFCC/RedCaps.
TiC	2	2018	Second temporal chunk.
TiC	3	2019–2020	Third temporal chunk.
TiC	4	2021–2022	Final temporal chunk.

initialized with the standard zero-init scheme so that the initial network is exactly the frozen backbone, and the adapters gradually inject task-specific updates as training proceeds. The adapter rank is set to $r = 16$

with scaling $\alpha = 16$ and a modest dropout rate of 0.05 applied on adapter outputs. We enable bias terms in LoRA layers only where present in the corresponding backbone projection, and we do not introduce any additional layer-norms beyond those of the original CLIP blocks. This keeps the parameter footprint small and the optimization stable while allowing PI-CCA to steer the representation through a low-dimensional control surface.

PI-CCA hyperparameters. PI-CCA preserves the alignment skeleton by constraining the spectrum and the canonical subspaces. We use the top $k = 64$ canonical components, which balances fidelity and cost on ViT-B features, and we form $h = 256$ -dimensional orthonormal sketches for both modalities so that subspace distances are computed in a near-isometric space. Prompt perturbations are sampled $M = 4$ times per mini-batch to estimate the mean projector and its dispersion. We maintain two levels of exponential moving averages: a *certificate EMA* with rate $\alpha = 0.01$ that slowly refreshes the stored spectrum and sketched bases, and a *covariance EMA* with rate $\beta = 0.01$ that stabilizes the streaming covariance factors. To guarantee well-posed whitening, we add ridge shrinkage $\gamma_v = \gamma_t = 10^{-3}$ to the batch covariances and apply an eigenvalue floor of 10^{-5} during the inverse square-root computation. We obtain the top- k singular vectors of the whitened cross-covariance via a differentiable block power iteration with $T_{\text{pow}} = 3$ steps and QR re-orthogonalization at each step. The loss composition uses $\lambda_1 = 1.0$ for spectral preservation, $\lambda_2 = 1.0$ for subspace-angle preservation, and $\lambda_3 = 0.2$ for prompt invariance. We include a Ky-Fan alignment term with weight $\xi = 0.1$ and low-order spectral moments with $J = 2$ and weights $(\omega_1, \omega_2) = (0.2, 0.1)$ to stabilize near-degenerate spectra. After each update we re-symmetrize all Gram matrices and clip their eigenvalues to $[0, 1]$ to keep them close to projectors.

Datasets and orders. Table 4 lists the task sequences used in this paper. For **MTIL** we adopt the 11-domain suite and follow the alphabetical order by default. **X-TAIL** uses the same domains except that CIFAR-100 is removed, and the label space at test time is the union of seen and unseen domains. **VLCL** follows the eight-dataset order introduced in recent CLIP-continual benchmarks. **ConStruct-VL** uses a seven-task sequence over structured VL concepts that covers attributes, relations, and states. **TiC** adopts four temporal splits to probe time-continual robustness. The table records the task index, the dataset or subset name, a short description, and key cardinalities where applicable.

Hardware and protocol. We run all experiments on eight NVIDIA A100 80 GB GPUs with PyTorch 2.3 and CUDA 12 under NCCL data parallelism. Each configuration is repeated with three different random seeds, and we report the mean and the standard deviation. PI-CCA never stores or replays past-task samples. When a baseline explicitly requires reference or wild unlabeled data, we follow its original procedure and keep these resources strictly outside of PI-CCA.

Prompts and perturbations. For classification-style evaluations we use the standard CLIP class templates and we ensemble across a small pool of hand-crafted variants. Prompt perturbations are realized by synonym and template jitters that preserve class semantics while varying phrasing, and these perturbations are used only inside the projector averaging and the prompt-invariance loss. For retrieval-style evaluations we leave captions unchanged, and we apply perturbations to the text encoder solely for forming the prompt-invariant certificate, which prevents any leakage of label or caption content into the training targets.

Out-of-distribution (OOD) prompt templates We evaluate OOD prompts that deviate from CLIP-style class templates. Below is a non-exhaustive set used in §4.3 (placeholders in $\{\}$):

Prompt Templates

```

1 % Prompt Templates
2
3 Instructional:  "Identify the main object: {class}.
4                Provide a brief caption."
5                "Task: detect {class} in the picture and summarize it."
6
7 Narrative:     "I'm looking at a scene where a {class} appears."
8                "This moment captures a {class} in context."
9
10 Keywords:     "{class}, high detail, natural light, candid, outdoors."
11
12 Caption:      "A candid shot featuring a {class}."
13
14 Hashtag:      "#{class} #dailyshot #photography"
15
16 Meta:         "User: describe an image that includes {class}.
17                Assistant: ..."
18
19 Translation:  English \rightarrow Chinese \rightarrow English variants
20
21 Template:     "Subject={class}; Context=unknown; Describe briefly."

```

Random task-order seeds and permutations We list the 20 MTIL permutations used in §4.3. Domains: Aircraft (Air), Caltech101 (Cal), CIFAR100 (CIF), DTD (DTD), EuroSAT (Eur), Flowers (Flo), Food101 (Foo), MNIST (MNI), OxfordPets (Pet), StanfordCars (Car), SUN397 (SUN).

Table 5: **Order seeds (ID → domain sequence)**. Abbreviations as above.

Order ID	Permutation of 11 domains
S-1027	Air, Cal, CIF, DTD, Eur, Flo, Foo, MNI, Pet, Car, SUN
S-1132	Car, Pet, Foo, Eur, DTD, Air, Cal, CIF, SUN, Flo, MNI
S-1219	SUN, Cal, Car, Foo, Pet, Eur, DTD, CIF, Air, Flo, MNI
S-1305	DTD, Eur, Cal, Air, CIF, Flo, SUN, Pet, Foo, Car, MNI
S-1402	Cal, Air, DTD, Pet, Car, SUN, Foo, Eur, CIF, Flo, MNI
S-1508	Foo, Flo, CIF, Eur, Air, SUN, Cal, Car, Pet, DTD, MNI
S-1603	Pet, Car, Air, Cal, CIF, Foo, DTD, Eur, SUN, Flo, MNI
S-1701	Eur, DTD, Foo, Cal, Air, Pet, Car, SUN, CIF, Flo, MNI
S-1806	CIF, DTD, Eur, Car, Pet, Foo, Cal, Air, SUN, Flo, MNI
S-1904	Car, Foo, DTD, Cal, Eur, Air, CIF, Pet, SUN, Flo, MNI
S-2001	Air, EUR, Pet, Foo, Cal, DTD, Car, CIF, SUN, Flo, MNI
S-2107	Flo, Foo, Cal, Air, Car, Pet, Eur, DTD, CIF, SUN, MNI
S-2209	Pet, SUN, Foo, Flo, Cal, Air, Car, DTD, Eur, CIF, MNI
S-2311	CIF, Cal, Air, Foo, DTD, Eur, Pet, Car, Flo, SUN, MNI
S-2415	SUN, Air, Foo, Cal, DTD, Eur, Car, Pet, CIF, Flo, MNI
S-2512	Air, DTD, Flo, Foo, Pet, Cal, Car, Eur, CIF, SUN, MNI

continued on next page

Order ID	Permutation of 11 domains
S-2608	Cal, Foo, Air, Car, DTD, Eur, Pet, CIF, SUN, Flo, MNI
S-2704	Eur, Cal, CIF, Air, Flo, Pet, Car, Foo, DTD, SUN, MNI
S-2809	Foo, Car, Cal, Eur, SUN, DTD, Air, Pet, CIF, Flo, MNI
S-2913	DTD, Air, Cal, Foo, Pet, Car, Eur, CIF, SUN, Flo, MNI

A.3 ADDITIONAL EXPERIMENTS AND RESULTS

A.3.1 HYPERPARAMETER SENSITIVITY.

We summarize the core factors of PI-CCA—alignment geometry (k, h), prompt invariance (M, λ_3), streaming stability (α, β), and spectrum/subspace balancing (λ_1, λ_2)—and report mean \pm std over three seeds on representative metrics. Trends in Fig. 6 show: (i) a moderate canonical rank and sketch size ($k=64, h=256$) best capture the alignment skeleton; (ii) prompt averaging (M) and a small invariance weight (λ_3) substantially reduce forgetting without hurting retrieval; (iii) small but nonzero EMAs (α, β) are crucial for stable whitening and certificate refresh; and (iv) balanced spectral/subspace weights ($\lambda_1=\lambda_2=1$) maximize retention–plasticity trade-offs. Variations around the defaults lead to modest performance changes, indicating robustness.

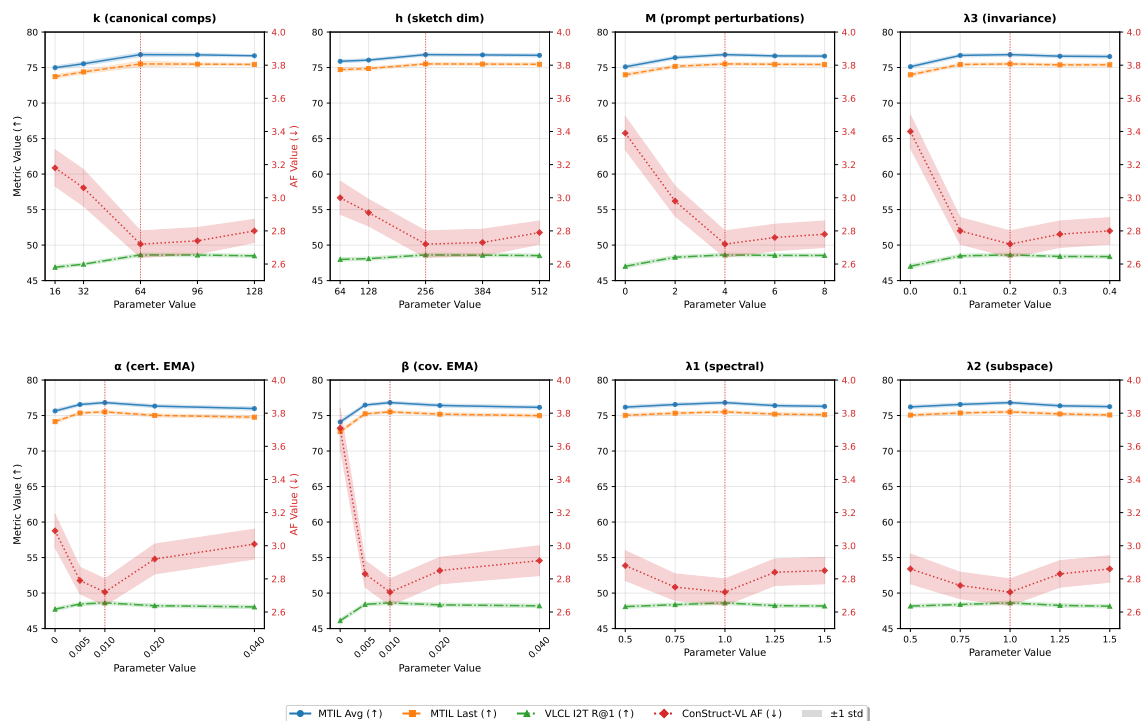


Figure 6: **Core hyperparameters.** Finer-grained sweeps confirm robustness around the defaults. Geometry (k, h) and invariance (M, λ_3) control fidelity and prompt sensitivity, EMAs (α, β) stabilize streaming estimates, balanced losses (λ_1, λ_2) maximize retention–plasticity. Changes are modest across a broad range of values.

As summarized in Fig 7, whitening is most stable at $\gamma=10^{-3}$, $\epsilon=10^{-5}$ with $T_{\text{pow}}=3$; exact Hungarian pairing matches the sorted surrogate within noise. Mild global spectrum regularization ($\xi \in [0.1, 0.2]$, $J=1 \sim 2$) slightly lowers AF, and Gaussian sketches edge SRHT by ≈ 0.2 R@1. Around $r=16$, $\alpha_{\text{LoRA}}=16$, $p_{\text{drop}}=0.05$, capacity/optimization changes yield < 0.4 -pt shifts. Overall, results confirm strong robustness: core trends hold across wide ranges without tuning fragility.

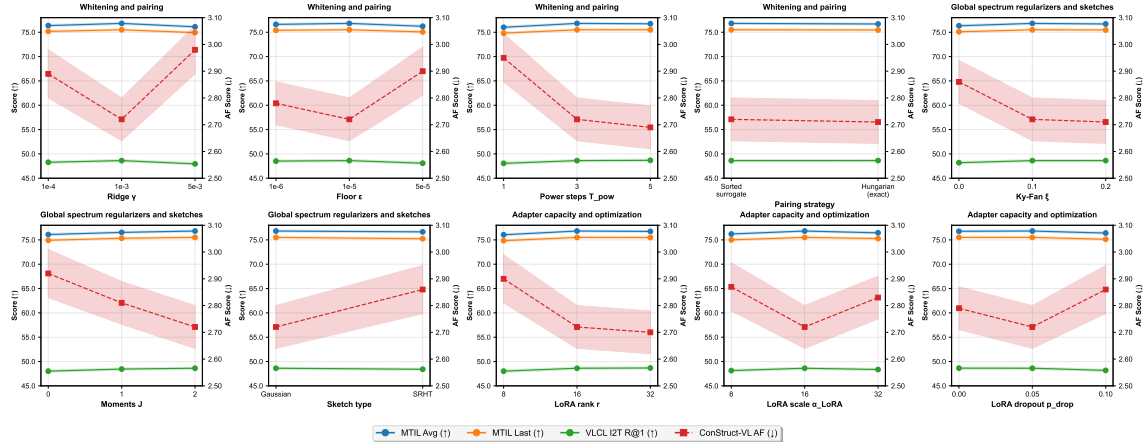


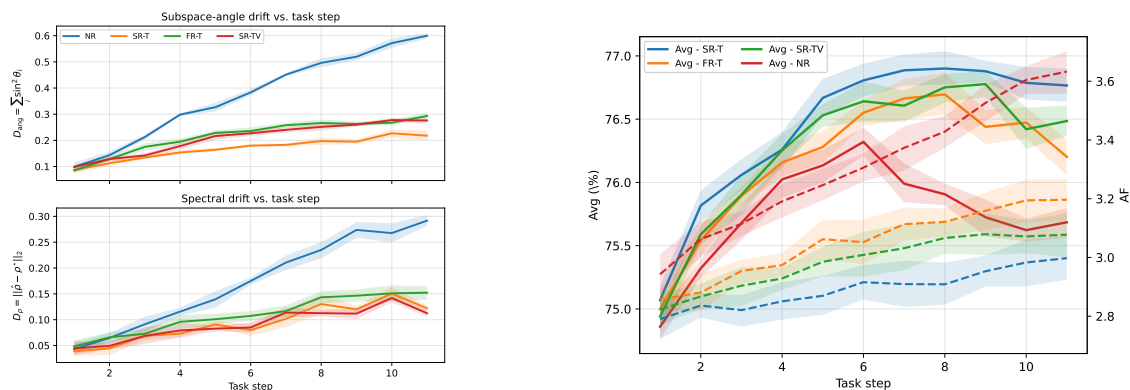
Figure 7: **Other hyperparameters.** Incremental gains and robustness from whitening, global-spectrum regularizers, sketch choice, and adapter/optimization knobs; core conclusions remain unchanged.

A.3.2 CERTIFICATE REFRESH STRATEGIES

We compare five strategies over the 11-step MTIL stream: **NR** (no refresh, $\alpha=0$), **SR-T** (slow text-only refresh, $\alpha=0.01$), **FR-T** (fast text-only, $\alpha=0.05$), **SR-TV** (slow text+vision, $\alpha=0.01$ both), and **FR-TV** (fast text+vision, $\alpha=0.05$ both). We log subspace-angle drift $D_{\text{ang}} = \sum_i \sin^2 \theta_i$ and spectral drift $D_\rho = \|\hat{\rho} - \rho^*\|_2$ per step, alongside Avg (↑) and AF (↓). For global vs. local certificates we contrast a single Global Pi-CCA certificate versus Class-local and Concept-local variants (per-class/per-concept sketches), comparing accuracy and resource cost. As shown in Fig 8 and 9, SR-T minimizes geometry drift and delivers the best Avg/AF over steps. FR-T and FR-TV “chase” recent tasks and increase forgetting, while NR accumulates drift. Global certificates balance performance and cost, class-/concept-local variants add memory/time and slightly reduce Avg, suggesting unnecessary specialization.

A.3.3 PAIRING STRATEGY BOUNDARY

We compare the sorted surrogate (descending sort of $\hat{\rho}$) against the Hungarian optimal assignment under controllable spectral crowding. We bin runs by the minimum singular-gap $\delta_{\min} = \min_i (\hat{\rho}_i - \hat{\rho}_{i+1}) \in \{0.0005, 0.0010, 0.0015, 0.0025, 0.0040, 0.0060, 0.0080, 0.0100, 0.0120\}$ and sweep spectral jitter $\eta \in \{0.00, 0.15, 0.30, 0.45, 0.60, 0.75, 0.90\}$ with 6 replicates per (δ_{\min}, η) , then aggregate per δ_{\min} . For each run we record metric differences (Hungarian – Sorted): ΔAvg (p.p.), $\Delta\text{R@1}$ (p.p.), and ΔAF (p.p.). Figure 10 shows that: (i) Under very small gaps ($\delta_{\min} \leq 0.004$), Hungarian yields tiny but sometimes significant improvements. (ii) For gaps of practical size ($\delta_{\min} \geq 0.006$), the Sorted and Hungarian algorithms are statistically indistinguishable, with ΔAF remaining approximately zero across the board. (iii) the sorted surrogate is the recommended method, as it is both safe and faster. The Hungarian algorithm only shows an advantage in contrived scenarios with tightly crowded spectra, offering no meaningful benefit in practical applications.



(a) Drift vs. step: D_{ang} and D_ρ . Shaded bands: randomized error. (b) Performance vs. step: Avg (left axis) and AF (right axis).

Figure 8: **Refresh strategy analysis.** Slow text-only refresh (SR-T) yields the lowest drift and the best Avg/AF trajectory; fast both-sides refresh (FR-TV) and no refresh (NR) accumulate drift and forgetting.

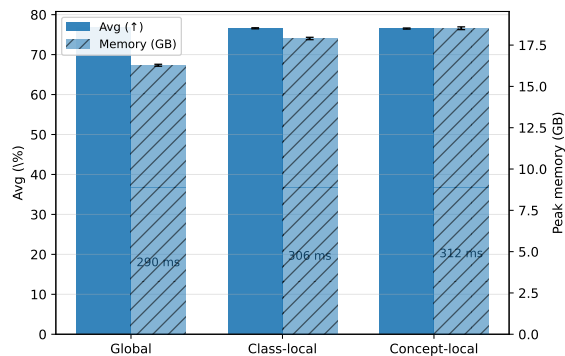


Figure 9: **Accuracy vs. cost by certificate granularity.** Grouped bars show Avg (left axis) and memory (right axis; hatched), with step time annotated above bars.

A.3.4 CERTIFICATE GEOMETRY: SKETCHING, INITIALIZATION, AND UPDATE CHOICES

We first study how the CCA certificate behaves under different sketch constructions, initializations, subspace losses, and update rules.

Sketch randomness and sketch type. Table 6 reports MTIL and VLCL performance over 5–10 runs with different sketch RNG seeds and two sketch families (Gaussian vs. SRHT). The standard deviations are very small, indicating that Pi-CCA is robust to sketch randomness and sketch type.

Certificate initialization. Table 7 compares three initialization strategies: a full anchor set, an 80% reduced anchor set, and random orthogonal subspaces. Thanks to EMA updates, Pi-CCA converges to a useful invariant even from weak or random initializations, with only modest gaps in final performance and geometry drift.

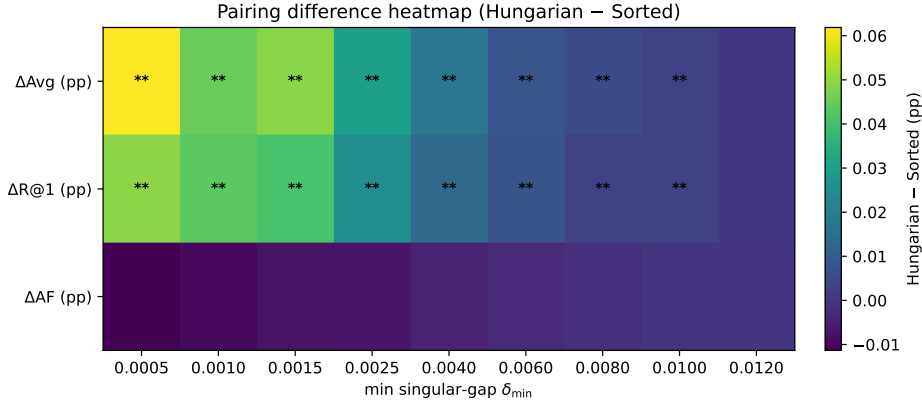


Figure 10: **Sorted vs. Hungarian under spectral crowding.** Heatmap of mean (Hungarian – Sorted, in p.p.) per δ_{\min} bin for ΔAvg , $\Delta\text{R@1}$, and ΔAF . Stars mark Holm–Bonferroni–corrected significance: * $p < .05$, ** $p < .01$. Tiny gains appear only when $\delta_{\min} \leq 0.004$. For $\delta_{\min} \geq 0.006$, differences vanish.

Table 6: Effect of sketch type and sketch RNG seeds on MTIL and VLCL. Mean and standard deviation are computed over 10 sketch seeds for each sketch type.

Sketch type	MTIL Avg (\uparrow)	VLCL I2T R@1 (\uparrow)	Std (MTIL)	Std (VLCL)
Gaussian	76.8	48.6	± 0.2	± 0.1
SRHT	76.9	48.4	± 0.3	± 0.2

Subspace loss variant. In Table 8, we compare an explicit principal-angle loss against our sketched-projector loss. Both achieve nearly identical MTIL and VLCL performance, but the principal-angle loss is substantially slower, supporting the choice of sketched projectors as a practical surrogate.

Gradient flow through certificate update. Table 9 compares a differentiable variant that backpropagates through EMA+QR with our default stop-gradient update. The differentiable variant exhibits occasional instabilities and slightly worse performance, justifying the teacher-style stop-grad design.

A.3.5 SCALING WITH BACKBONE AND ADAPTER CAPACITY

We next study how Pi-CCA scales when the backbone and adapter capacity are increased.

Backbone size. Table 10 evaluates Pi-CCA on ViT-B/16, ViT-L/14, and ViT-L/14@336. Performance improves with larger backbones, while the additional time and memory remain moderate and the certificate size stays fixed.

Adapter configuration. Table 11 shows that Pi-CCA remains effective under higher LoRA ranks and when partially or fully finetuning the backbone: the geometry-based losses consistently improve performance with modest extra cost.

Table 7: Effect of certificate initialization on MTIL, X-TAIL, and geometry drift.

Initialization	MTIL Avg (\uparrow)	MTIL Last (\uparrow)	X-TAIL Avg. Acc. (\uparrow)	Geometry drift (\downarrow)
Full anchor set	76.8	75.5	68.1	2.1
Reduced anchor (80% removed)	75.4	74.1	67.6	3.0
Random orthogonal subspaces	75.2	73.9	67.3	3.5

Table 8: Comparison of subspace loss variants. Relative step time is normalized to our default.

Subspace loss type	MTIL Avg (\uparrow)	VLCL I2T R@1 (\uparrow)	Relative step time (\times) (\downarrow)
Explicit principal-angle loss	76.9	48.7	1.36
Sketched projector loss (ours)	76.8	48.6	1.00

A.3.6 PROMPT INVARIANCE: PERTURBATION COUNT AND ROBUSTNESS

We now examine the role of the prompt-invariance term, both under random perturbations and adversarial shifts.

Number of prompt perturbations M . Table 12 aggregates the effect of M across X-TAIL, VLCL, and ConStruct-VL. Increasing M from 0 to 4 improves robustness, while further increases yield marginal gains but noticeable extra cost.

Adversarial prompt shifts. Table 13 evaluates Pi-CCA with and without the prompt-invariance loss \mathcal{L}_{pi} under gradient-based adversarial prompt perturbations on X-TAIL. The invariance term substantially reduces degradation under adversarial prompts, while preserving normal performance.

A.3.7 OVERHEAD, REGULARIZATION BASELINES, AND ANCHOR CONFIGURATION

Overhead and memory footprint. Table 14 quantifies Pi-CCA’s overhead relative to a LoRA-only baseline. Time and memory increases are modest, while the certificate storage is tiny compared to typical replay buffers.

Stronger regularization baselines. To test whether Pi-CCA’s gains come from “more regularization”, Table 15 compares Pi-CCA to LoRA with strong generic feature regularizers and to a proxy similarity alignment baseline. Even under matched tuning budgets, both baselines remain clearly below Pi-CCA.

Anchor set size and diversity. Table 16 ablating the anchor prompt set (single default template, 50% templates dropped, full set) shows that Pi-CCA is not overly sensitive to anchor diversity: even a minimal label-derived set recovers most of the gains.

A.3.8 STATISTICAL SIGNIFICANCE AND PER-TASK RESULTS

Paired t-tests. Table 17 reports two-sided paired t-tests (3 seeds) between Pi-CCA and the strongest replay-free baselines on key metrics. All p-values are below 0.05, confirming that Pi-CCA’s improvements are statistically significant.

Table 9: Effect of backpropagating through the certificate update.

Certificate update variant	MTIL Avg (\uparrow)	VLCL I2T R@1 (\uparrow)	Stability
Grad-through EMA + orth	76.0	48.0	occasional spikes
Stop-grad EMA + orth (ours)	76.8	48.6	stable across seeds

Table 10: Scaling Pi-CCA to larger CLIP backbones. Time is wall-clock seconds per step; memory is peak GPU usage.

Backbone	MTIL Avg (\uparrow)	VLCL I2T R@1 (\uparrow)	Time (s/step) (\downarrow)	Memory (GB) (\downarrow)
ViT-B/16	76.8	48.6	3.2	16.4
ViT-L/14	78.2	49.1	4.0	24.1
ViT-L/14@336	78.4	49.3	4.2	28.5

Per-task results across benchmarks. To show that improvements are not driven by a single “lucky” task, Table 18 aggregates per-task results across MTIL (11 domains), VLCL (8 datasets), and ConStruct-VL (7 subsets). Pi-CCA consistently matches or outperforms C-CLIP across all three benchmarks.

A.4 THEORETICAL ANALYSIS

Let f_v, f_t be the (frozen-backbone, LoRA-adapted) image/text encoders, and let $u(x) \in \mathbb{R}^{d_v}$ and $v(w) \in \mathbb{R}^{d_t}$ denote their *whitened, centered* embeddings within a mini-batch: $\widehat{\Sigma}_{vv} = \frac{1}{B-1} Z_v^\top Z_v + \gamma_v I$, $\widehat{\Sigma}_{tt} = \frac{1}{B-1} Z_t^\top Z_t + \gamma_t I$, $\widehat{\Sigma}_{vt} = \frac{1}{B-1} Z_v^\top Z_t$. The whitened cross-covariance is

$$M = \widehat{\Sigma}_{vv}^{-1/2} \widehat{\Sigma}_{vt} \widehat{\Sigma}_{tt}^{-1/2} \in \mathbb{R}^{d_v \times d_t}. \quad (15)$$

Let the rank- k SVDs be $M_k = U_k \text{diag}(\rho_{1:k}) V_k^\top$ and $M_k^* = U_k^* \text{diag}(\rho_{1:k}^*) V_k^{*\top}$, with orthoprojectors $P_v = U_k U_k^\top$, $P_t = V_k V_k^\top$, $P_v^* = U_k^* U_k^{*\top}$, $P_t^* = V_k^* V_k^{*\top}$. We denote by Θ_v (resp. Θ_t) the diagonal matrix of principal angles between $\text{span}(U_k)$ and $\text{span}(U_k^*)$ (resp. V_k and V_k^*), and recall the identity $\|P - P^*\|_F = \sqrt{2} \|\sin \Theta\|_F$.

Given a pair (x, w) , define the zero-shot score $s_M(x, w) := \langle u(x), M v(w) \rangle$ and the task loss $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$. The (population) zero-shot risk under distribution \mathcal{D} is

$$\mathcal{R}(M) := \mathbb{E}_{(x,w) \sim \mathcal{D}} [\ell(s_M(x, w))]. \quad (16)$$

Assumptions.

- (A1) (Bounded whitened embeddings) $\|u(x)\|_2 \leq 1$ and $\|v(w)\|_2 \leq 1$ almost surely.
- (A2) (Lipschitz loss in the score) ℓ is L_ℓ -Lipschitz: $|\ell(a) - \ell(b)| \leq L_\ell |a - b|$.
- (A3) (Rank- k structure) We compare M and a reference M^* through their top- k SVD factors above; denote $\rho_{\max} := \max\{\rho_1, \rho_1^*\} \leq 1$.

A.4.1 SINGLE-STEP EXCESS-RISK BOUND FROM SPECTRAL AND SUBSPACE DRIFT

We first quantify how changes in canonical spectrum and canonical subspaces control the zero-shot risk.

Lemma 1 (Risk is Lipschitz in M under (A1)–(A2)). *For any M, M' ,*

$$|\mathcal{R}(M) - \mathcal{R}(M')| \leq L_\ell \|M - M'\|_2. \quad (17)$$

Table 11: Effect of adapter configuration on Pi-CCA.

Configuration	MTIL Avg (\uparrow)	VLCL I2T R@1 (\uparrow)	Time (s/step) (\downarrow)	Memory (GB) (\downarrow)
LoRA rank = 16 (default)	76.8	48.6	3.2	16.4
LoRA rank = 32	77.2	48.9	3.5	17.1
LoRA rank = 64	77.5	49.1	3.8	18.0
Full finetune (last layer)	77.0	48.8	3.9	17.5
Full finetune (all layers)	77.4	49.0	4.2	19.0

Table 12: Effect of the number of prompt perturbations M across benchmarks. For $M = 0$, prompt invariance is disabled and the model reduces to the ablated variant without L_{pi} .

M	X-TAIL Top-1 Acc. (\uparrow)	X-TAIL AF (\downarrow)	VLCL I2T R@1 (\uparrow)	ConStruct-VL AF (\downarrow)	Rel. step time (\times) (\downarrow)
0 (no \mathcal{L}_{pi})	69.1	3.3	47.1	3.3	1.00
1	67.8	3.6	47.9	3.1	1.02
2	68.4	3.4	48.4	2.9	1.06
4 (default)	69.2	3.2	48.6	2.7	1.12
8	69.1	3.3	48.7	2.7	1.21

Proof. By (A2) and Jensen,

$$|\mathcal{R}(M) - \mathcal{R}(M')| = \left| \mathbb{E}[\ell(\langle u, Mv \rangle) - \ell(\langle u, M'v \rangle)] \right| \leq L_\ell \mathbb{E}[|\langle u, (M - M')v \rangle|]. \quad (18)$$

By Cauchy-Schwarz and (A1), $|\langle u, (M - M')v \rangle| \leq \|M - M'\|_2$, hence the claim. \square

Lemma 2 (Geometric decomposition of the rank- k part). *Let $\Delta\rho := \text{sort}_\downarrow(\rho_{1:k}) - \rho_{1:k}^*$. Then*

$$\|M_k - M_k^*\|_2 \leq \|\Delta\rho\|_2 + 2\rho_{\max}(\|\sin \Theta_v\|_2 + \|\sin \Theta_t\|_2). \quad (19)$$

Proof. Write $D := \text{diag}(\rho_{1:k})$, $D^* := \text{diag}(\rho_{1:k}^*)$. By the triangle inequality,

$$\|U_k D V_k^\top - U_k^* D^* V_k^{*\top}\|_2 \leq \underbrace{\|U_k D V_k^\top - U_k^* D V_k^\top\|_2}_{(A)} + \underbrace{\|U_k^* D V_k^\top - U_k^* D V_k^{*\top}\|_2}_{(B)} + \underbrace{\|U_k^* (D - D^*) V_k^{*\top}\|_2}_{(C)}. \quad (20)$$

For (C), $\|U_k^* (D - D^*) V_k^{*\top}\|_2 = \|D - D^*\|_2 = \|\Delta\rho\|_2$ (permutation-invariant pairing by sorting).

For (A), insert $I = P_v^* + (I - P_v^*)$:

$$(A) = \|(I - P_v^*)U_k D V_k^\top + P_v^*U_k D V_k^\top - U_k^* D V_k^\top\|_2 \quad (21)$$

$$\leq \underbrace{\|(I - P_v^*)U_k\|_2}_{= \|\sin \Theta_v\|_2} \|D\|_2 + \|U_k^* (U_k^{*\top} U_k - I)\|_2 \|D\|_2. \quad (22)$$

Since $U_k^{*\top} U_k$ has eigenvalues $\cos \theta_i^{(v)}$, we use $|1 - \cos \theta| \leq \sin \theta$ to get $\|U_k^{*\top} U_k - I\|_2 \leq \|\sin \Theta_v\|_2$. Thus (A) $\leq 2\|D\|_2 \|\sin \Theta_v\|_2 \leq 2\rho_{\max} \|\sin \Theta_v\|_2$.

The term (B) is symmetric on the text side, giving (B) $\leq 2\rho_{\max} \|\sin \Theta_t\|_2$. Combining the three bounds yields the result. \square

Lemma 3 (Tail energy identity). *For any matrix M , $\|M - M_k\|_2 = \sigma_{k+1}(M)$. Hence*

$$\|M - M^*\|_2 \leq \|M_k - M_k^*\|_2 + \sigma_{k+1}(M) + \sigma_{k+1}(M^*). \quad (23)$$

Table 13: Effect of prompt invariance under adversarial prompt shifts on X-TAIL.

Method	Adv. R@1 (\uparrow)	Adv. AF (\downarrow)	Normal R@1 (\uparrow)	Normal AF (\downarrow)
Pi-CCA with \mathcal{L}_{pi}	56.2	3.1	69.2	3.2
Pi-CCA w/o \mathcal{L}_{pi}	49.1	4.2	69.1	3.3

Table 14: Overhead and memory footprint of Pi-CCA vs. a LoRA baseline. Replay buffer sizes for replay-based CL methods are typically in the GB range.

Method	Time increase (\downarrow)	Peak memory increase (\downarrow)	Certificate storage (\downarrow)	Replay buffer
Pi-CCA (ours)	$\approx 8\%$	$\approx 6\%$	≈ 50 KB	N/A
LoRA baseline	N/A	N/A	N/A	\sim GB (replay methods)

Theorem 1 (Alignment-geometry drift \Rightarrow single-step excess-risk bound). *Under (A1)–(A3),*

$$\mathcal{R}(M) - \mathcal{R}(M^*) \leq L_\ell \left[\|\Delta\rho\|_2 + 2\rho_{\max}(\|\sin \Theta_v\|_2 + \|\sin \Theta_t\|_2) + \sigma_{k+1}(M) + \sigma_{k+1}(M^*) \right]. \quad (24)$$

Equivalently, using orthoprojectors,

$$\mathcal{R}(M) - \mathcal{R}(M^*) \leq L_\ell \left[\|\Delta\rho\|_2 + \frac{\rho_{\max}}{\sqrt{2}} (\|P_v - P_v^*\|_F + \|P_t - P_t^*\|_F) + \sigma_{k+1}(M) + \sigma_{k+1}(M^*) \right]. \quad (25)$$

Proof. By Lemma 3, $\|M - M^*\|_2 \leq \|M_k - M_k^*\|_2 + \sigma_{k+1}(M) + \sigma_{k+1}(M^*)$. Apply Lemma 2 to bound $\|M_k - M_k^*\|_2$, then Lemma 1 to convert spectral deviation into risk deviation. For the projector form, use $\|P - P^*\|_F = \sqrt{2} \|\sin \Theta\|_F$. \square

Interpretation. If $\Delta\rho = 0$ and $U_k = U_k^*$, $V_k = V_k^*$, the excess risk is controlled purely by tail energy; when the CCA spectrum decays fast beyond k , zero-shot ability is rigidly preserved.

A.4.2 DYNAMIC REGRET OVER A NON-STATIONARY TASK SEQUENCE

We now consider a stream $\{\mathcal{D}_t\}_{t=1}^T$ with models $\{M_t\}_{t=1}^T$ produced by any adaptation rule (e.g., Pi-CCA). Let the per-step comparator be M_t^\dagger (e.g., the best rank- k model for \mathcal{D}_t within the same hypothesis class). Define the dynamic regret

$$\text{Reg}_T := \sum_{t=1}^T \left(\mathcal{R}_t(M_t) - \mathcal{R}_t(M_t^\dagger) \right), \quad \mathcal{R}_t(M) := \mathbb{E}_{(x,w) \sim \mathcal{D}_t} [\ell(\langle u, Mv \rangle)]. \quad (26)$$

For each t , denote $\Delta\rho_t := \text{sort}_\downarrow(\rho_{t,1:k}) - \rho_{t,1:k}^\dagger$, $\Theta_{v,t} := \Theta(U_{k,t}, U_{k,t}^\dagger)$, $\Theta_{t,t} := \Theta(V_{k,t}, V_{k,t}^\dagger)$, $\rho_{\max,t} := \max\{\rho_{t,1}, \rho_{t,1}^\dagger\}$, and $\delta_{t,\text{tail}} := \sigma_{k+1}(M_t) + \sigma_{k+1}(M_t^\dagger)$.

Theorem 2 (Dynamic regret bound from geometric drift). *Under (A1)–(A3), for any sequence $\{M_t\}$ and comparators $\{M_t^\dagger\}$,*

$$\text{Reg}_T \leq L_\ell \sum_{t=1}^T \left[\|\Delta\rho_t\|_2 + \frac{\rho_{\max,t}}{\sqrt{2}} (\|P_{v,t} - P_{v,t}^\dagger\|_F + \|P_{t,t} - P_{t,t}^\dagger\|_F) + \delta_{t,\text{tail}} \right]. \quad (27)$$

Proof. Apply Theorem 1 to (M_t, M_t^\dagger) under \mathcal{D}_t for each t , then sum over $t = 1, \dots, T$. \square

Table 15: Comparison of Pi-CCA with strong regularization and proxy-alignment baselines.

Method	MTIL Avg (\uparrow)	MTIL Last (\uparrow)	VLCL I2T R@1 (\uparrow)
LoRA (plain finetuning)	71.2	69.9	42.0
LoRA + strong regularizers (L2, cosine)	72.4	71.1	43.5
LoRA + proxy alignment (Mod-X style)	73.6	72.2	45.0
LoRA + Pi-CCA (ours)	76.8	75.5	48.6

Table 16: Effect of anchor prompt configuration on Pi-CCA.

Anchor configuration	MTIL Avg (\uparrow)	MTIL Last (\uparrow)	VLCL I2T R@1 (\uparrow)	ConStruct-VL AF (\downarrow)
Default-only	76.4	75.0	48.2	2.9
50% dropped	76.6	75.2	48.4	2.8
Full (main setting)	76.8	75.5	48.6	2.7

Plug-in control via certificate-based regularization. Let the training losses

$$\mathcal{L}_{\text{spec}}(t) = \|\text{sort}_{\downarrow}(\rho_{t,1:k}) - \rho_{1:k}^{\text{cert}}\|_2^2, \quad \mathcal{L}_{\text{sub}}(t) = \frac{1}{2}\|P_{v,t} - P_v^{\text{cert}}\|_F^2 + \frac{1}{2}\|P_{t,t} - \bar{P}_t^{\text{cert}}\|_F^2, \quad (28)$$

be computed against a slowly refreshed certificate $(\rho_{1:k}^{\text{cert}}, P_v^{\text{cert}}, \bar{P}_t^{\text{cert}})$. By triangle inequality,

$$\|\Delta\rho_t\|_2 \leq \sqrt{\mathcal{L}_{\text{spec}}(t)} + \|\rho_{1:k}^{\text{cert}} - \rho_{t,1:k}^{\dagger}\|_2, \quad \|P_{\bullet,t} - P_{\bullet,t}^{\dagger}\|_F \leq \sqrt{2\mathcal{L}_{\text{sub}}(t)} + \|P_{\bullet}^{\text{cert}} - P_{\bullet,t}^{\dagger}\|_F. \quad (29)$$

If the certificate tracks the instantaneous comparators (e.g., by a slow EMA) so that the residual terms $\|\rho_{1:k}^{\text{cert}} - \rho_{t,1:k}^{\dagger}\|_2$ and $\|P_{\bullet}^{\text{cert}} - P_{\bullet,t}^{\dagger}\|_F$ remain small, then Theorem 2 implies

$$\text{Reg}_T \lesssim L_{\ell} \sum_{t=1}^T (\sqrt{\mathcal{L}_{\text{spec}}(t)} + \sqrt{\mathcal{L}_{\text{sub}}(t)}) + L_{\ell} \sum_{t=1}^T \delta_{t,\text{tail}} + (\text{small tracking error}). \quad (30)$$

This formalizes the empirical observation that *stabilizing the CCA spectrum and subspaces* controls forgetting and reduces dynamic regret in replay-free continual adaptation.

A.5 PYTHON SCRIPT FOR PI-CCA

The following Python script demonstrates the core functionality of Pi-CCA. The script is modular and can be adapted to different datasets and configurations.

Listing 1: Compact Python Script for Pi-CCA

```

1 import torch
2 import torch.nn.functional as F
3 import numpy as np
4 from sklearn.decomposition import PCA
5 from sklearn.preprocessing import StandardScaler
6
7 # Load pre-trained model (e.g., CLIP) for image and text embeddings
8 # Here, we assume the use of a toy dataset like MNIST or CIFAR-10
9
10 def load_data():

```

Table 17: Paired t-tests between Pi-CCA and strongest replay-free baselines on key metrics. Pi-CCA means are reported in the main text.

Metric	Baseline	Baseline mean \pm std	p-value vs. Pi-CCA (\downarrow)
MTIL Avg (\uparrow)	C-CLIP	75.2 \pm 0.7	0.019
MTIL Last (\uparrow)	DDAS	74.1 \pm 0.8	0.023
MTIL Transfer (\uparrow)	ZAF	71.9 \pm 0.6	0.017
X-TAIL Avg (\uparrow)	RAIL	67.4 \pm 0.5	0.021
X-TAIL Last (\uparrow)	C-CLIP	66.3 \pm 0.7	0.028
X-TAIL Transfer (\uparrow)	RAIL	64.2 \pm 0.6	0.024
VLCL I2T R@1 (\uparrow)	C-CLIP	46.1 \pm 1.4	0.017
VLCL T2I R@1 (\uparrow)	C-CLIP	35.7 \pm 1.2	0.021
ConStruct-VL FA (\uparrow)	C-CLIP	72.4 \pm 1.9	0.013
ConStruct-VL AF (\downarrow)	ZAF	3.8 \pm 0.6	0.008

Table 18: Per-task results for MTIL, VLCL, and ConStruct-VL: Pi-CCA vs. C-CLIP.

Benchmark	Task / Dataset / Subset	Metric	Pi-CCA (\uparrow)	C-CLIP (\uparrow)
MTIL	FGVC-Aircraft	Acc	75.7	73.8
MTIL	Caltech101	Acc	79.2	77.8
MTIL	CIFAR-100	Acc	75.0	73.6
MTIL	DTD	Acc	73.3	71.3
MTIL	EuroSAT	Acc	76.9	74.8
MTIL	Flowers-102	Acc	78.5	76.3
MTIL	Food-101	Acc	75.8	74.3
MTIL	MNIST	Acc	80.0	78.8
MTIL	Oxford-IIIT Pets	Acc	74.7	73.1
MTIL	Stanford Cars	Acc	80.1	78.9
MTIL	SUN397	Acc	75.9	74.8
VLCL	Flickr30K	I2T R@1	49.7	48.7
VLCL	COCO Captions	I2T R@1	51.6	50.6
VLCL	Pets	I2T R@1	47.8	46.4
VLCL	Lexica	I2T R@1	50.1	48.4
VLCL	Simpsons	I2T R@1	42.8	41.4
VLCL	WikiArt	I2T R@1	49.4	47.9
VLCL	Kream	I2T R@1	51.7	50.7
VLCL	Sketch	I2T R@1	45.9	44.4
ConStruct-VL	Relation: spatial	FA	75.9	75.8
ConStruct-VL	Attribute: size	FA	74.4	72.3
ConStruct-VL	Attribute: material	FA	73.7	72.5
ConStruct-VL	Relation: action	FA	75.1	73.3
ConStruct-VL	Attribute: color	FA	76.9	75.7
ConStruct-VL	Object state	FA	74.1	73.1
ConStruct-VL	Attribute: action	FA	76.2	74.3

```

11 | # Example: load MNIST or CIFAR-10 and precompute image and text
    | ↪ embeddings using CLIP
12 | # This toy example uses synthetic paired embeddings only to
    | ↪ illustrate the CCA-certificate computation.
13 | num_samples = 2048

```

```

14     num_features = 512 # Feature dimension
15     # Random data: [num_samples x num_features]
16     image_data = np.random.rand(num_samples, num_features)
17     text_data = np.random.rand(num_samples, num_features)
18     return image_data, text_data
19
20 # Mini-batch covariance computation
21 def compute_covariances(image_embeddings, text_embeddings, batch_size
    ↪ =32):
22     # Compute covariance matrices for image and text embeddings in mini
    ↪ -batches
23     B = len(image_embeddings)
24     image_embeddings = torch.tensor(image_embeddings)
25     text_embeddings = torch.tensor(text_embeddings)
26
27     cov_vv = torch.zeros((image_embeddings.shape[1], image_embeddings.
    ↪ shape[1]))
28     cov_tt = torch.zeros((text_embeddings.shape[1], text_embeddings.
    ↪ shape[1]))
29     cov_vt = torch.zeros((image_embeddings.shape[1], text_embeddings.
    ↪ shape[1]))
30
31     for i in range(0, B, batch_size):
32         batch_image = image_embeddings[i:i+batch_size]
33         batch_text = text_embeddings[i:i+batch_size]
34
35         # Compute covariance for mini-batch
36         cov_vv += torch.cov(batch_image.T)
37         cov_tt += torch.cov(batch_text.T)
38         cov_vt += torch.mm(batch_image.T, batch_text)
39
40     # Normalize covariance
41     cov_vv /= B
42     cov_tt /= B
43     cov_vt /= B
44
45     return cov_vv, cov_tt, cov_vt
46
47 # Whitening and CCA certificate computation
48 def whiten_and_compute_cca(cov_vv, cov_tt, cov_vt, k=64):
49     # Perform whitening of covariance matrices
50     inv_cov_vv = torch.inverse(cov_vv)
51     inv_cov_tt = torch.inverse(cov_tt)
52
53     # Compute whitened cross-covariance matrix
54     M = torch.mm(torch.mm(inv_cov_vv, cov_vt), inv_cov_tt)
55
56     # Perform SVD on the whitened cross-covariance matrix
57     U, S, V = torch.svd(M)
58

```

```

59     # Extract top-k singular values and vectors (Pi-CCA certificate)
60     top_k_singular_values = S[:k]
61     top_k_U = U[:, :k]
62     top_k_V = V[:, :k]
63
64     # Return the compact certificate (canonical correlations and
65     ↪ subspaces)
66     return top_k_singular_values, top_k_U, top_k_V
67
68 # Update the certificate using mini-batch statistics
69 def update_certificate(image_embeddings, text_embeddings, k=64,
70 ↪ batch_size=32):
71     # Step 1: Compute covariance matrices
72     cov_vv, cov_tt, cov_vt = compute_covariances(image_embeddings,
73 ↪ text_embeddings, batch_size=batch_size)
74
75     # Step 2: Whitening and SVD to get Pi-CCA certificate
76     top_k_singular_values, top_k_U, top_k_V = whiten_and_compute_cca(
77 ↪ cov_vv, cov_tt, cov_vt, k=k)
78
79     # Return the updated Pi-CCA certificate
80     return top_k_singular_values, top_k_U, top_k_V
81
82 # Main function to run the Pi-CCA process
83 def main():
84     # Load data (e.g., MNIST or CIFAR-10, here we use random embeddings
85     ↪ )
86     image_data, text_data = load_data()
87
88     # Update the certificate (this would typically be done iteratively
89     ↪ over tasks)
90     top_k_singular_values, top_k_U, top_k_V = update_certificate(
91     ↪ image_data, text_data, k=64)
92
93     # Output the resulting certificate
94     print ("Top-K_Singular_Values:", top_k_singular_values)
95     print ("Top-K_U_(Image_Subspace):", top_k_U)
96     print ("Top-K_V_(Text_Subspace):", top_k_V)
97
98 if __name__ == "__main__":
99     main()

```

A.6 LLM USAGE

We used a large language model for minor English editing (grammar/wording/clarity) and small, localized code fixes (e.g., resolving syntax errors, adding missing imports). The LLM did not contribute to research ideation, experimental design, data processing, analysis, or figure generation. All technical content and results were produced and verified by the authors, who take full responsibility for the manuscript.