

CNSP: CONSISTENT NULL-SPACE PROJECTION FOR PRINCIPLED PROMPT-BASED CONTINUAL LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Continual learning aims to acquire new knowledge sequentially without forgetting previous tasks, yet catastrophic forgetting remains a major challenge. Prompt-based continual learning has recently shown competitive empirical progress, yet its theoretical underpinnings remain incomplete. We introduce Consistent Null-Space Projection (CNSP), the first unified and mathematically rigorous framework for representational consistency in prompt-based continual learning. CNSP proves that task-performance preservation reduces to two jointly sufficient requirements—feature preservation and head preservation—while deriving explicit consistency conditions under full Transformer parameterization. These conditions yield a tractable null-space projection rule for stable prompt updates. Across various benchmarks and backbones, CNSP demonstrates consistent improvements in accuracy and forgetting, with especially clear benefits in high-dimensional and domain-shift scenarios.

1 INTRODUCTION

Continual learning seeks to enable models to acquire new knowledge from sequential tasks without forgetting previously learned knowledge. Despite significant progress, catastrophic forgetting (McCloskey & Cohen, 1989), where performance on old tasks deteriorates when learning new ones, remains a major challenge (Wang et al., 2024).

Recently, prompt-tuning methods, particularly Visual Prompt Tuning (VPT) (Jia et al., 2022) with Vision Transformers (ViTs) (Dosovitskiy et al., 2020), have emerged as a promising rehearsal-free solution for continual learning (Wang et al., 2024). By updating only lightweight prompts while freezing the backbone, these methods achieve both computational and storage efficiency. While several prompt-based methods for continual learning (e.g., L2P (Wang et al., 2022b), DualPrompt (Wang et al., 2022a), Coda-Prompt (Smith et al., 2023)) have demonstrated competitive empirical performance, theoretical foundations of prompt-based continual learning remain underdeveloped. Classical orthogonal-projection methods (Saha et al., 2021; Deng et al., 2021; Wang et al., 2021) guarantee retention by deriving explicit consistency conditions in linear or convolutional networks, but these guarantees do not transfer to Transformers: nonlinear LayerNorm, learned QKV projections, and multi-head attention invalidate linear assumptions.

To bridge this theoretical gap, recent efforts have adapted projection ideas to ViTs/VPT (Qiao et al., 2024; Lu et al., 2024). PGP (Qiao et al., 2024) approximates Multi-Head Self-Attention (MHSA) and LayerNorm as linear mappings and constructs orthogonal subspaces in the sum space of inputs and prompts, but its strong simplification limits generality. NSP² (Lu et al., 2024) advances this line of work by deriving approximate sufficient conditions for preserving representations under prompt updates. However, NSP² still leaves several fundamental aspects of the Transformer pipeline theoretically unresolved: the extension from single-head to multi-head attention lacks principled justification; LayerNorm’s broadcasted affine structure is simplified in ways that obscure its algebraic properties; strong assumptions on prompt distribution (mean–variance invariance) can introduce instability; and the classification head is omitted from the overall consistency formulation. Consequently, existing analyses offer valuable local insights but do not yield an end-to-end account of how prompt updates interact with the entire Transformer computation to preserve task performance.

Motivated by the absence of a complete theoretical closure, we revisit prompt-based continual learning from an end-to-end, matrix-level perspective and introduce Consistent Null-Space Projection (CNSP)—a principled framework that establishes the first unified representational-consistency formulation for prompted Transformers in continual learning. Our central insight is that performance preservation admits a clean mathematical decomposition into two jointly sufficient components: *feature preservation*, ensuring that MHSA+LayerNorm representations for previous tasks remain stable, and head preservation, guaranteeing decision-level consistency for previously learned classifier(s). Building on this decomposition, we derive sufficient consistency conditions for preserving past-task performance by jointly accounting for the coupled roles of MHSA, LayerNorm, and classification heads under full prompted Transformer parameterization, from which a tractable null-space projection rule emerges as an analytical consequence.

Our contributions are threefold:

(i) **Unified theoretical closure.** We establish the first unified and mathematically rigorous framework demonstrating that prompt-based continual learning reduces to a single consistency principle integrating feature preservation and head preservation.

(ii) **Sufficient consistency conditions.** By rigorously modeling MHSA, LayerNorm, and classifier behavior using an algebraically precise matrix-form characterization under full prompted Transformer parameterization, we derive explicit sufficient conditions for past-task performance preservation, leading to a tractable null-space projection rule.

(iii) **Consistent empirical gains.** CNSP reliably reduces forgetting and improves accuracy across diverse benchmarks and backbones, with more noticeable gains in high-dimensional and cross-domain settings—all without introducing architectural heuristics or task-specific prompt designs.

2 RELATED WORKS

2.1 PROMPT-BASED METHODS FOR CONTINUAL LEARNING

Prompt-based methods have emerged as a promising paradigm for continual learning, offering lightweight task adaptation with reduced interference. Early approaches include L2P (Wang et al., 2022b), which employs a learnable prompt pool with dynamic prompt selection, and DualPrompt (Wang et al., 2022a), which further disentangles prompts into general and expert subsets to balance knowledge sharing and task-specific adaptation. CPrompt (Gao et al., 2024) addresses train–test inconsistency via classifier- and prompt-consistent learning, while CODA-Prompt (Smith et al., 2023) decomposes prompts into fine-grained components and dynamically recombining them through attention. EvoPrompt (Kurniawan et al., 2024) formulates prompt learning as evolutionary search for long-term adaptability. Beyond purely visual prompts, LGCL (Khan et al., 2023) incorporates semantic guidance by aligning visual representations with language embeddings at both task and class levels. ConvPrompt (Roy et al., 2024) leverages convolutional structures to construct hierarchical prompts and employs large language models to estimate task similarity.

More recent efforts pursue principled formulations. PGP (Qiao et al., 2024) enforces gradient orthogonality in the joint input–prompt space by linearizing the nonlinear effects of MHSA and LayerNorm, while NSP² (Lu et al., 2024) derives null-space projection conditions for VPT by explicitly analyzing attention and LayerNorm. CPG (Lu et al., 2025) extends this direction with consistency-constrained mixtures of experts.

Prompt-based continual learning methods are versatile and effective, but their theoretical foundations remain limited. NSP² advances this line via null-space projection for VPT, yet its simplified treatment of multi-head attention and LayerNorm reduces robustness and interpretability. We address these gaps by formalizing MHSA and LayerNorm in exact matrix form and deriving sufficient conditions for performance preservation, thus yielding a principled and effective framework.

2.2 ORTHOGONAL PROJECTION METHODS FOR CONTINUAL LEARNING

Orthogonal projection methods have been widely studied in CNNs and MLPs (e.g., OWM (Zeng et al., 2019), GPM (Saha et al., 2021), Adam-NSCL (Wang et al., 2021)) as a principled way to reduce task interference. The core idea is to preserve past feature responses by constrain-

ing parameter updates within subspaces orthogonal to previous tasks’ features, formalized as $\mathbf{X}^\top(\Theta + \Delta\Theta) = \mathbf{X}^\top\Theta$, where \mathbf{X} denotes input features, Θ the convolutional or linear parameters, and $\Delta\Theta$ the update. Enforcing $\mathbf{X}^\top\Delta\Theta = \mathbf{0}$ guarantees invariance of past intermediate representations, thereby retaining knowledge. This principle, rooted in gradient limitation theory, provides a mathematical explanation of the stability-plasticity trade-off. Representative methods include GPM Saha et al. (2021), which projects gradients onto the null space of previously learned inputs; TRGP (Lin et al., 2022b), which refines GPM by scaling task gradients with a trust region matrix before projection; and Connector (Lin et al., 2022a) interpolates a normally updated model with one constrained by projection.

These methods work well in linear or convolutional architectures, where orthogonality condition holds exactly. However, in Transformers, nonlinear MHSA, softmax, and LayerNorm interactions violate this linearity, making classical projections unreliable. Recent works have applied orthogonal projection to prompt tuning (Qiao et al., 2024; Lu et al., 2024). NSP² (Lu et al., 2024) unfolds VPT forward propagation and derives approximate sufficient conditions for performance preservation via null-space projection. Building on this foundation, we introduce stricter and more general consistency conditions alongside stable optimization strategies, thereby strengthening projection-based continual learning in VPT.

3 PRELIMINARIES

In this section, we introduce the notations and conventions, describe the forward propagation process of transformer blocks under VPT-deep (Jia et al., 2022)—where learnable prompts are inserted into the input token sequence at each block for efficient adaptation with a frozen backbone—and formalize the continual learning problem in this setting.

3.1 NOTATIONS

We use the following notational conventions: (i) Non-bold letters denote positive integers, e.g., $a, A \in \mathbb{N}^+$. (ii) Bold lowercase denote vectors, $\mathbf{a} \in \mathbb{R}^n$; bold uppercase denote matrices, $\mathbf{A} \in \mathbb{R}^{n \times m}$. (iii) Column concatenation: for $\mathbf{A}^{(i)} \in \mathbb{R}^{n \times d}$, $[\mathbf{A}^{(i)}]_{i=1}^N = [\mathbf{A}^{(1)} \dots \mathbf{A}^{(N)}] \in \mathbb{R}^{n \times Nd}$. (iv) Row concatenation: for $\mathbf{A}^{(i)} \in \mathbb{R}^{n \times d}$, $[\mathbf{A}^{(1)}; \dots; \mathbf{A}^{(N)}] \in \mathbb{R}^{Nn \times d}$. (v) Slicing: if $\mathbf{A} = [\mathbf{A}^{(1)}; \mathbf{A}^{(2)}]$ with $\mathbf{A}^{(1)} \in \mathbb{R}^{n \times d}$, then $\mathbf{A}[:n] = \mathbf{A}^{(1)}$.

3.2 FORWARD PROPAGATION IN VPT-DEEP

Let $\mathbf{X} \in \mathbb{R}^{(N+1) \times D}$ denote the patch embeddings of an input image, where N is the number of patches, the first row is the pre-trained [CLS] token, and D the embedding dimension. A standard ViT with depth L can be written as $f_{\text{ViT}}(\mathbf{X} \mid \Theta; \{\mathbf{W}, \mathbf{b}\})$, where Θ are frozen backbone parameters and $\{\mathbf{W}, \mathbf{b}\}$ are learnable classification head parameters. Let $\mathbf{X}^{(l)} \in \mathbb{R}^{(N+1) \times D}$ denote the input embeddings of the l -th block.

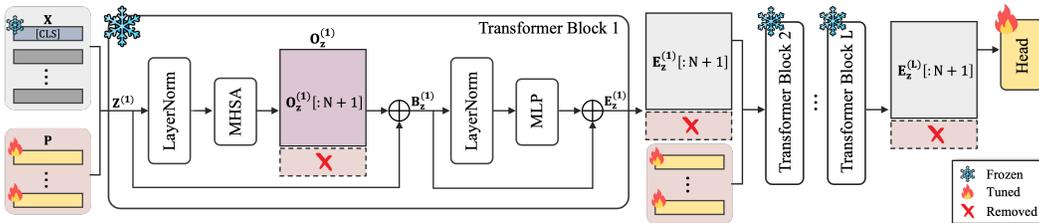


Figure 1: Forward propagation in VPT. Image embeddings (\mathbf{X}) and learnable prompt embeddings (\mathbf{P}) are jointly processed by LayerNorm, MHSA, and MLP with residual connections in each transformer block. Modules marked with ❄️ are frozen, those with 🔥 are trainable, and those with ❌ indicate information discarded before entering the next layer.

In VPT-deep, each block l introduces M learnable prompts $\mathbf{P}^{(l)} \in \mathbb{R}^{M \times D}$, concatenated with image embeddings to form the input sequences: $\mathbf{Z}^{(l)} = [\mathbf{X}^{(l)}; \mathbf{P}^{(l)}] \in \mathbb{R}^{(N+1+M) \times D}$, $\mathbf{X}^{(1)} = \mathbf{X}$. The forward propagation in VPT-deep is illustrated in fig. 1. Each transformer block applies LayerNorm (LN), Multi-Head Self-Attention (MHSA), and an MLP with residual connections:

$$\mathbf{O}_Z^{(l)} = \text{MHSA}(\text{LN}(\mathbf{Z}^{(l)})), \quad \mathbf{B}_Z^{(l)} = \mathbf{Z}^{(l)} + \mathbf{O}_Z^{(l)}, \quad \mathbf{E}_Z^{(l)} = \mathbf{B}_Z^{(l)} + \text{MLP}(\text{LN}(\mathbf{B}_Z^{(l)})). \quad (1)$$

Before the next block, the prompts are removed and replaced with a new set of learnable prompts:

$$\mathbf{X}^{(l+1)} = \mathbf{E}_Z^{(l)}[: N + 1], \quad \mathbf{Z}^{(l+1)} = [\mathbf{X}^{(l+1)}; \mathbf{P}^{(l+1)}]. \quad (2)$$

Thus, prompts modulate intermediate features at each block but do not persist across blocks. The final [CLS] token is passed into the classification head:

$$f_{\text{ViT}}(\mathbf{X} \mid \Theta; \{\mathbf{W}, \mathbf{b}\}; \{\mathbf{P}^{(l)}\}_{l=1}^L) = \text{softmax}(\mathbf{E}_Z^{(L)}[: 1]\mathbf{W} + \mathbf{b}), \quad (3)$$

where $\mathbf{W} \in \mathbb{R}^{D \times C}$, $\mathbf{b} \in \mathbb{R}^C$, and C is the number of classes.

3.3 PROBLEM FORMULATION: CONTINUAL LEARNING WITH VPT-DEEP

We consider a continual learning setting with a sequence of tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$. Each task \mathcal{T}_t is associated with a dataset $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{N_t}$, where $\mathbf{x}_t^{(i)}$ is an image and $y_t^{(i)}$ its label. At block l , let $\mathbf{P}_t^{(l)}$ be the prompts after learning tasks up to \mathcal{T}_t . For the next task \mathcal{T}_{t+1} , updates are:

$$\mathbf{P}_{t+1}^{(l)} = \mathbf{P}_t^{(l)} + \Delta \mathbf{P}^{(l)}, \quad \mathbf{W}_{t+1} = \mathbf{W}_t + \Delta \mathbf{W}, \quad \mathbf{b}_{t+1} = \mathbf{b}_t + \Delta \mathbf{b}, \quad (4)$$

The continual learning objective requires preserving performance on previous tasks after training on a new one. We formalize the objective as:

$$f_{\text{ViT}}(\mathbf{X}_t \mid \Theta; \{\mathbf{W}_{t+1}, \mathbf{b}_{t+1}\}; \{\mathbf{P}_{t+1}^{(l)}\}_{l=1}^L) = f_{\text{ViT}}(\mathbf{X}_t \mid \Theta; \{\mathbf{W}_t, \mathbf{b}_t\}; \{\mathbf{P}_t^{(l)}\}_{l=1}^L), \quad (5)$$

which enforces that model outputs on old inputs \mathbf{X}_t from \mathcal{T}_t must remain unchanged. Equivalently, updates $\{\Delta \mathbf{P}^{(l)}\}_{l=1}^L$ and $\{\Delta \mathbf{W}, \Delta \mathbf{b}\}$ must not affect the predictions for previously learned tasks. To characterize this objective, we introduce two key propositions that decompose it into two complementary requirements: **feature preservation** at intermediate blocks and **head preservation** at the classification head.

Proposition 1. Consider two tasks \mathcal{T}_t and \mathcal{T}_{t+1} . If the image-token part of attention outputs remains unchanged for all blocks in VPT, i.e.,

$$\mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(l)}[: N + 1] = \mathbf{O}_{\mathbf{Z}_{t,t}}^{(l)}[: N + 1], \quad l = 1, 2, \dots, L, \quad (6)$$

then the image-token output of the final block is preserved:

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[: N + 1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[: N + 1], \quad (7)$$

which we refer to as **feature preservation**. Here $\mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(l)}$, $\mathbf{O}_{\mathbf{Z}_{t,t}}^{(l)}$, $\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}$, and $\mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}$ are computed by eq. (1). In particular, the final [CLS] token is preserved:

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[: 1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[: 1]. \quad (8)$$

Proof Sketch. By eq. (2), since prompts are replaced at each block, their effect propagates only through the image-token outputs $\mathbf{E}_Z^{(l)}[: N + 1]$. Row-wise independence of post-MHSA operations ensures invariance of $\mathbf{O}_Z^{(l)}[: N + 1]$ carries forward. Full proof is provided in Appendix A.

Proposition 2. Consider tasks \mathcal{T}_t and \mathcal{T}_{t+1} . If the final [CLS] token remains unchanged (i.e., eq. (8) holds) and the head updates satisfy:

$$\mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[: 1]\Delta \mathbf{W} + \Delta \mathbf{b} = \mathbf{0}, \quad (9)$$

which we refer to as **head preservation**, then the continual learning objective eq. (5) is satisfied.

216 *Proof.* From eq. (8) and eq. (9), we have

$$217 \mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[:1]\mathbf{W}_{t+1} + \mathbf{b}_{t+1} = \mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[:1](\mathbf{W}_t + \Delta\mathbf{W}) + (\mathbf{b}_t + \Delta\mathbf{b}) = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[:1]\mathbf{W}_t + \mathbf{b}_t, \quad (10)$$

219 which, by eq. (3), implies eq. (5). \square

221 Together, Propositions 1 and 2 reduce continual learning preservation to two conditions: feature
222 preservation, requiring prompt updates to leave image-token representations invariant for represen-
223 tational stability, and head preservation, requiring the classifier to maintain consistent mappings of
224 the preserved [CLS] token. This decomposition of objective eq. (5) forms the theoretical basis for
225 the sufficient conditions derived in the next section. Particularly, proposition 1 implies interme-
226 diate feature preservation reduces to ensuring invariance at the MHSA output. Therefore, in the
227 derivations presented in the next section, we focus on LayerNorm and MHSA, as they are the only
228 operations directly relevant to feature preservation.

229 4 METHOD

231 Our method builds on Propositions 1 and 2, enforcing both feature and head preservation. Feature
232 preservation is ensured through explicit constraints on prompt updates and head preservation guar-
233 anteed by design. According to proposition 1, the preservation objective can be analyzed blockwise
234 by focusing solely on the MHSA and its preceding LayerNorm operation. Other components (e.g.,
235 MLP) are row-wise and do not alter the preservation property. Since eq. (6) must be satisfied at
236 every layer, we omit the superscript (l) when deriving algebraic sufficient conditions on prompt
237 updates.

239 4.1 UNFOLDING FORWARD PROPAGATION OF LAYERNORM AND MHSA IN A SINGLE 240 BLOCK

241 Let the input sequence be $\mathbf{Z}_{t,t+1} = [\mathbf{X}_t; \mathbf{P}_{t+1}] \in \mathbb{R}^{(N+1+M) \times D}$, where $\mathbf{P}_{t+1} = \mathbf{P}_t + \Delta\mathbf{P}$. Figure 2
242 illustrates the unrolled forward pass of LayerNorm (LN) and MHSA in a single transformer block.
243 In this setting, computations involving $\mathbf{Q}_{\mathbf{P}_{t+1}}^{(h)}$ can be safely omitted; justification is provided below.

244 First, $\mathbf{Z}_{t,t+1}$ is passed through LayerNorm, we introduce the following lemma:

245 **Lemma 1.** *Consider the computation of LayerNorm $\text{LN}(\cdot)$. Let $\mathbf{A} = [\mathbf{A}^{(1)}; \mathbf{A}^{(2)}]$ with $\mathbf{A}^{(1)} \in$
246 $\mathbb{R}^{n \times d}$ and $\mathbf{A}^{(2)} \in \mathbb{R}^{m \times d}$. Then $\text{LN}(\mathbf{A}) = [\text{LN}(\mathbf{A}^{(1)}); \text{LN}(\mathbf{A}^{(2)})]$.*

247 Proof of Lemma 1 is in Appendix A. According to lemma 1, we have $\text{LN}(\mathbf{Z}_{t,t+1}) =$
248 $[\text{LN}(\mathbf{X}_t); \text{LN}(\mathbf{P}_t + \Delta\mathbf{P})]$. Next, the normalized sequence undergoes QKV transformation, in-
249 cluding QKV projections and multi-head splitting. For head $h \in \{1, \dots, H\}$, we have:

$$250 \Phi_{\mathbf{Z}_{t,t+1}}^{(h)} = \text{LN}(\mathbf{Z}_{t,t+1}) \mathbf{W}_\phi^{(h)} + \mathbf{1}_{N+1+M} \mathbf{b}_\phi^{(h)\top} = \begin{bmatrix} \text{LN}(\mathbf{X}_t) \mathbf{W}_\phi^{(h)} + \mathbf{1}_{N+1} \mathbf{b}_\phi^{(h)\top} \\ \text{LN}(\mathbf{P}_t + \Delta\mathbf{P}) \mathbf{W}_\phi^{(h)} + \mathbf{1}_M \mathbf{b}_\phi^{(h)\top} \end{bmatrix}, \quad (11)$$

251 where $\Phi \in \{\mathbf{Q}, \mathbf{K}, \mathbf{V}\}$, ϕ serves as the symbolic index for q, k, v , and $\mathbf{W}_\phi^{(h)} \in \mathbb{R}^{D \times D_H}$, $\mathbf{b}_\phi^{(h)} \in$
252 \mathbb{R}^{D_H} , and $D_H = D/H$. Let

$$253 \mathbf{Q}_{\mathbf{Z}_{t,t+1}}^{(h)} = \begin{bmatrix} \mathbf{Q}_{\mathbf{X}_t}^{(h)} \\ \mathbf{Q}_{\mathbf{P}_{t+1}}^{(h)} \end{bmatrix}, \quad \mathbf{K}_{\mathbf{Z}_{t,t+1}}^{(h)} = \begin{bmatrix} \mathbf{K}_{\mathbf{X}_t}^{(h)} \\ \mathbf{K}_{\mathbf{P}_{t+1}}^{(h)} \end{bmatrix}, \quad \mathbf{V}_{\mathbf{Z}_{t,t+1}}^{(h)} = \begin{bmatrix} \mathbf{V}_{\mathbf{X}_t}^{(h)} \\ \mathbf{V}_{\mathbf{P}_{t+1}}^{(h)} \end{bmatrix}, \quad (12)$$

254 where $\mathbf{Q}_{\mathbf{X}_t}^{(h)}, \mathbf{K}_{\mathbf{X}_t}^{(h)}, \mathbf{V}_{\mathbf{X}_t}^{(h)} \in \mathbb{R}^{(N+1) \times D_H}$ and $\mathbf{Q}_{\mathbf{P}_{t+1}}^{(h)}, \mathbf{K}_{\mathbf{P}_{t+1}}^{(h)}, \mathbf{V}_{\mathbf{P}_{t+1}}^{(h)} \in \mathbb{R}^{M \times D_H}$.

255 For each head h , the attention weights and value aggregation are computed as

$$256 \mathbf{H}_{\mathbf{Z}_{t,t+1}}^{(h)} = \text{softmax}\left(\mathbf{Q}_{\mathbf{Z}_{t,t+1}}^{(h)} \mathbf{K}_{\mathbf{Z}_{t,t+1}}^{(h)\top} / \sqrt{D_H}\right) \mathbf{V}_{\mathbf{Z}_{t,t+1}}^{(h)}, \quad (13)$$

257 where the softmax function acts on the rows of matrix $\mathbf{Q}_{\mathbf{Z}_{t,t+1}}^{(h)} \mathbf{K}_{\mathbf{Z}_{t,t+1}}^{(h)\top}$. Finally, concatenating
258 across heads and applying the output projection gives

$$259 \mathbf{O}_{\mathbf{Z}_{t,t+1}} = \left[\mathbf{H}_{\mathbf{Z}_{t,t+1}}^{(h)}\right]_{h=1}^H \mathbf{W}_o + \mathbf{1}_{N+1+M} \mathbf{b}_o^\top. \quad (14)$$

By proposition 1, feature preservation requires $\mathbf{O}_{Z_t, t+1}[: N + 1] = \mathbf{O}_{Z_t, t}[: N + 1]$. Note that

$$[\mathbf{H}_{Z_t, t+1}^{(h)}]_{h=1}^H[: N + 1] = \left[\text{softmax}(\mathbf{Q}_{X_t}^{(h)} \mathbf{K}_{Z_t, t+1}^{(h)\top} / \sqrt{D_H}) \mathbf{V}_{Z_t, t+1}^{(h)} \right]_{h=1}^H. \quad (15)$$

Therefore, only $\mathbf{Q}_{X_t}^{(h)}$ contributes to the image-token outputs, terms from $\mathbf{Q}_{P_{t+1}}^{(h)}$ can be omitted in the preservation analysis.

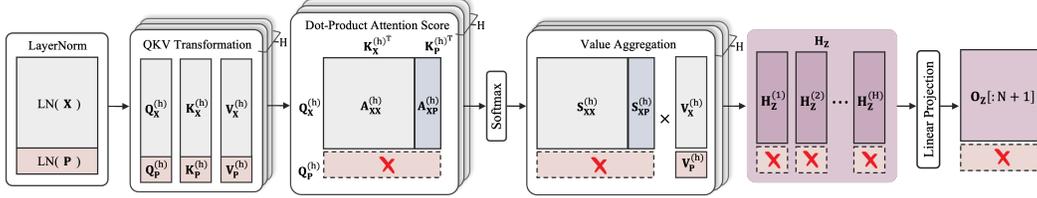


Figure 2: Detailed illustration of LayerNorm and the multi-head self-attention mechanism in a VPT transformer block, showing how image and prompt embeddings are jointly processed. Those with \times indicate information discarded before entering the next layer.

4.2 SUFFICIENT CONDITIONS FOR FEATURE PRESERVATION IN VPT

Building on the above expansion, we restrict attention to the $\mathbf{Q}_{X_t}^{(h)}$ -dependent pathways. Define

$$\mathbf{A}_{Z_t, t+1}^{(h)} = \frac{\mathbf{Q}_{X_t}^{(h)} \mathbf{K}_{Z_t, t+1}^{(h)\top}}{\sqrt{D_H}} = \frac{1}{\sqrt{D_H}} \mathbf{Q}_{X_t}^{(h)} \begin{bmatrix} \mathbf{K}_{X_t}^{(h)\top} & \mathbf{K}_{P_{t+1}}^{(h)\top} \end{bmatrix} \triangleq \begin{bmatrix} \mathbf{A}_{X_t X_t}^{(h)} & \mathbf{A}_{X_t P_{t+1}}^{(h)} \end{bmatrix}, \quad (16)$$

$$\mathbf{S}_{Z_t, t+1}^{(h)} = \text{softmax} \left(\begin{bmatrix} \mathbf{A}_{X_t X_t}^{(h)} & \mathbf{A}_{X_t P_{t+1}}^{(h)} \end{bmatrix} \right) \triangleq \begin{bmatrix} \mathbf{S}_{X_t X_t}^{(h)} & \mathbf{S}_{X_t P_{t+1}}^{(h)} \end{bmatrix}, \quad (17)$$

where $\mathbf{A}_{X_t X_t}^{(h)}, \mathbf{S}_{X_t X_t}^{(h)} \in \mathbb{R}^{(N+1) \times (N+1)}$ and $\mathbf{A}_{X_t P_{t+1}}^{(h)}, \mathbf{S}_{X_t P_{t+1}}^{(h)} \in \mathbb{R}^{(N+1) \times M}$. (Symbol-module correspondence is depicted in 2.) By proposition 1, eq. (14) and eq. (15), feature preservation requires

$$\left[\mathbf{S}_{X_t X_t}^{(h)} \mathbf{V}_{X_t}^{(h)} + \mathbf{S}_{X_t P_{t+1}}^{(h)} \mathbf{V}_{P_{t+1}}^{(h)} \right]_{h=1}^H \mathbf{W}_o = \left[\mathbf{S}_{X_t X_t}^{(h)} \mathbf{V}_{X_t}^{(h)} + \mathbf{S}_{X_t P_t}^{(h)} \mathbf{V}_{P_t}^{(h)} \right]_{h=1}^H \mathbf{W}_o \quad (18)$$

In MHSA, \mathbf{W}_o fuses per-head outputs into a unified representation (Vaswani et al., 2017; Horn & Johnson, 2012). Since \mathbf{W}_o may in principle be singular, we impose stronger sufficient conditions by requiring per-head equality:

$$\mathbf{S}_{X_t P_{t+1}}^{(h)} \mathbf{V}_{P_{t+1}}^{(h)} - \mathbf{S}_{X_t P_t}^{(h)} \mathbf{V}_{P_t}^{(h)} = \mathbf{0}. \quad (19)$$

Direct algebraic characterization is intractable due to the nonlinearity of softmax. Following NSP² (Lu et al., 2024), we instead require invariance of the attention scores:

$$\mathbf{A}_{Z_t, t}^{(h)} = \mathbf{A}_{Z_t, t+1}^{(h)} \Leftrightarrow \mathbf{Q}_{X_t}^{(h)} \left(\mathbf{K}_{P_{t+1}}^{(h)\top} - \mathbf{K}_{P_t}^{(h)\top} \right) = \mathbf{0} \Rightarrow \mathbf{S}_{Z_t, t}^{(h)} = \mathbf{S}_{Z_t, t+1}^{(h)}. \quad (20)$$

That is, the attention scores from image tokens to prompts remain unchanged before and after the prompt update. Under eq. (20), condition eq. (19) reduces to

$$\mathbf{S}_{X_t P_t}^{(h)} \left(\mathbf{V}_{P_{t+1}}^{(h)} - \mathbf{V}_{P_t}^{(h)} \right) = \mathbf{0}. \quad (21)$$

By substituting eq. (11), we obtain per-head sufficient conditions for feature preservation:

$$\begin{cases} \mathbf{Q}_{X_t}^{(h)} \left(\mathbf{K}_{P_{t+1}}^{(h)\top} - \mathbf{K}_{P_t}^{(h)\top} \right) = \mathbf{Q}_{X_t}^{(h)} \mathbf{W}_k^{(h)\top} (\text{LN}(P_t + \Delta P) - \text{LN}(P_t))^\top = \mathbf{0} \\ \mathbf{S}_{X_t P_t}^{(h)} \left(\mathbf{V}_{P_{t+1}}^{(h)} - \mathbf{V}_{P_t}^{(h)} \right) = \mathbf{S}_{X_t P_t}^{(h)} (\text{LN}(P_t + \Delta P) - \text{LN}(P_t)) \mathbf{W}_v^{(h)} = \mathbf{0} \end{cases} \quad (22)$$

To reduce the conditions to direct dependence on ΔP , we expand the LayerNorm:

$$\text{LN}(\mathbf{P}_t + \Delta P) = \frac{\mathbf{P}_t + \Delta P - \boldsymbol{\mu}_{\mathbf{P}_t + \Delta P}}{\boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P}} \cdot \boldsymbol{\gamma} + \boldsymbol{\beta}, \quad (23)$$

where $\boldsymbol{\mu}_{\mathbf{P}_t + \Delta P}, \boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P} \in \mathbb{R}^M$ are row-wise statistics and $\boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^D$ are shared affine parameters (broadcast in implementation). To enable an algebraically exact characterization of prompt interactions, we explicitly reformulate LayerNorm in matrix form. A detailed derivation is provided in Appendix B. Let $\mathbf{C} = \mathbf{I}_D - \frac{1}{D} \mathbf{1}_D \mathbf{1}_D^\top (\in \mathbb{R}^{D \times D})$ be the centering matrix, $\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P}} (\in \mathbb{R}^{M \times M})$ ¹ the diagonal scaling matrix from $\boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P}$, and $\boldsymbol{\Gamma} (\in \mathbb{R}^{D \times D})$ the diagonal matrix from $\boldsymbol{\gamma}$. Then

$$\text{LN}(\mathbf{P}_t + \Delta P) = \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P}}^{-1} (\mathbf{P}_t + \Delta P) \mathbf{C} \boldsymbol{\Gamma} + \mathbf{1}_M \boldsymbol{\beta}^\top. \quad (24)$$

Note that pretrained $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are frozen. If we introduce a variance-invariance assumption

$$\boldsymbol{\sigma}_{\mathbf{P}_t} = \boldsymbol{\sigma}_{\mathbf{P}_t + \Delta P} \triangleq \boldsymbol{\sigma}, \quad (25)$$

then

$$\text{LN}(\mathbf{P}_t + \Delta P) - \text{LN}(\mathbf{P}_t) = \mathbf{D}_{\boldsymbol{\sigma}}^{-1} \Delta P \mathbf{C} \boldsymbol{\Gamma}. \quad (26)$$

Substituting eq. (26) into eq. (22), the sufficient constraints reduce to the following linear form:

$$\begin{cases} \mathbf{Q}_{\mathbf{X}_t}^{(h)} \mathbf{W}_k^{(h)\top} \boldsymbol{\Gamma} \mathbf{C} \Delta P^\top = \mathbf{0} \\ \mathbf{S}_{\mathbf{X}_t \mathbf{P}_t}^{(h)} \mathbf{D}_{\boldsymbol{\sigma}}^{-1} \Delta P \mathbf{C} \boldsymbol{\Gamma} \mathbf{W}_v^{(h)} = \mathbf{0} \end{cases} \quad h = 1, 2, \dots, H. \quad (27)$$

The first constraint ensures that, after centering and scaling, prompt updates do not perturb the scoring pattern from image tokens to prompts. The second ensures that they do not perturb the value aggregation pathway. Therefore, under the assumption eq. (25), eq. (27) provides a set of per-head sufficient conditions for feature preservation eq. (7).

4.3 OPTIMIZATION UNDER FEATURE PRESERVATION CONSTRAINTS

To enforce the variance-invariance assumption eq. (25), we adopt a soft constraint by introducing a standard deviation alignment loss:

$$\mathcal{L} = \mathcal{L}_{\text{ce}} + \lambda \mathcal{L}_{\text{std}}, \quad \mathcal{L}_{\text{std}} = \|\boldsymbol{\sigma}_{\mathbf{P}_{t+1}} - \boldsymbol{\sigma}_{\mathbf{P}_t}\|_1, \quad (28)$$

where \mathcal{L} is the total loss, \mathcal{L}_{ce} is the cross-entropy loss, and λ balances the alignment term. Under this assumption, we obtain the sufficient conditions in eq. (27), constraining ΔP . We adopt a right-side nullification form on the second constraint $\mathbf{S}_{\mathbf{X}_t \mathbf{P}_t}^{(h)} \mathbf{D}_{\boldsymbol{\sigma}}^{-1} \Delta P \mathbf{C} \boldsymbol{\Gamma} \mathbf{W}_v^{(h)} = \mathbf{0}$. The rationale is to avoid circular dependence: if $\mathbf{D}_{\boldsymbol{\sigma}}^{-1}$ remains on the solving side, then ΔP becomes entangled with the current statistic $\boldsymbol{\sigma}$, complicating implementation and hurting numerical stability. In contrast, the right-sided formulation $\Delta P \mathbf{C} \boldsymbol{\Gamma} \mathbf{W}_v^{(h)} = \mathbf{0}$ removes explicit dependence on $\boldsymbol{\sigma}$, yielding more stable behavior. Moreover, it only involves frozen pretrained parameters, reducing the computational overhead. Let $\mathbf{R}_k^{(h)} = \mathbf{C} \boldsymbol{\Gamma} \mathbf{W}_k^{(h)} \mathbf{Q}_{\mathbf{X}_t}^{(h)\top}$, $\mathbf{R}_v^{(h)} = \mathbf{C} \boldsymbol{\Gamma} \mathbf{W}_v^{(h)}$, and define

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_k^{(1)} & \dots & \mathbf{R}_k^{(H)} & \mathbf{R}_v^{(1)} & \dots & \mathbf{R}_v^{(H)} \end{bmatrix} \in \mathbb{R}^{D \times 2D}. \quad (29)$$

Our implementation enforces, at every optimization step, $\Delta P \mathbf{R} = \mathbf{0}$. Equivalently, each row of ΔP lies in the left null space of \mathbf{R} . Since the effective constraint size grows with tasks, handling \mathbf{R} directly can be expensive. We thus use the following basic result:

Lemma 2. For $\mathbf{A} \in \mathbb{R}^{n \times m}$ with left null space $\mathcal{N}_L(\mathbf{A}) = \{\mathbf{x} \mathbf{A} = \mathbf{0} \mid \mathbf{x}^\top \in \mathbb{R}^n\}$, it holds that $\mathcal{N}_L(\mathbf{A}) = \mathcal{N}_L(\mathbf{A} \mathbf{A}^\top)$.

Proof of lemma 2 is in Appendix A. Lemma 2 implies that $\Delta P \mathbf{R} = \mathbf{0}$ is equivalent to requiring each row of ΔP to lie in $\mathcal{N}_L(\mathbf{R} \mathbf{R}^\top)$, where $\mathbf{R} \mathbf{R}^\top \in \mathbb{R}^{D \times D}$ has fixed dimension during training.

¹In practice, a small constant $\epsilon > 0$ is added to the variance when computing each row’s standard deviation, ensuring positivity and numerical stability. Hence all standard deviations are strictly greater than zero, and the corresponding diagonal matrix is invertible.

In practice, we perform Singular Value Decomposition (SVD) on $\mathbf{R}\mathbf{R}^\top$, collect the left singular vectors corresponding to zero singular values to form an orthogonal basis \mathbf{U}_0 , so that $\mathcal{N}_L(\mathbf{R}\mathbf{R}^\top) = \text{span}(\mathbf{U}_0)$ and its null-space projector is $\mathbf{B}_R = \mathbf{U}_0\mathbf{U}_0^\top$. Let $\Delta\mathbf{P}_{\text{raw}}$ denote the raw update from back-propagation. We apply the feature-preservation projection

$$\Delta\mathbf{P} \leftarrow \Delta\mathbf{P}_{\text{raw}}\mathbf{B}_R, \quad (30)$$

which guarantees that the feature-preservation sufficient conditions (i.e., eq. (27)) are satisfied.

4.4 CLASSIFICATION HEAD PRESERVATION

Proposition 2 requires the classification head to maintain its mapping of the [CLS] token across tasks. This condition is naturally fulfilled by task-specific heads, where only the current task’s head is updated and all previous heads are frozen. Given CNSP’s feature-preservation constraint, the representations fed to earlier heads remain invariant, so their input–output mappings are preserved automatically. At inference, we concatenate logits from all task-specific heads and predict the class with highest confidence. An illustration is provided in fig. 4, with further discussion in Appendix F.1.

Summary. CNSP enforces feature preservation through null-space projected updates and ensures head preservation through frozen task-specific classifiers, resulting in end-to-end representational and decision-level consistency across tasks. Detailed algorithmic descriptions are provided in Appendix C, and comparison with NSP² is included in Appendix D.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETTINGS

Benchmarks and Evaluation Metrics. We evaluate CNSP on three benchmarks widely used in prompt-based continual learning: 10/20-Split CIFAR-100 (Krizhevsky et al., 2009), 10-Split ImageNet-R (Hendrycks et al., 2021), and 10-Split DomainNet (Peng et al., 2019). We also include a domain-incremental benchmark, 6-Split-DomainNet, where each task contains a single distinct visual domain. Performance is measured by Last Average Accuracy (ACC, \uparrow) and Forgetting (\downarrow). Details of the benchmarks and metric definitions are provided in Appendix E.1 and E.2.

Implementation. Our main experiments use three backbones: ImageNet-21K pretrained ViT-B/16 (Dosovitskiy et al., 2020), DINOv2 ViT-B/14 (Oquab et al., 2023), and AugReg ViT-L/16 (Steiner et al., 2021), each equipped with four learnable prompt tokens per layer. Experiments follow the class-incremental setting with disjoint classes and unknown task identity at inference, and are trained with Adam (Kingma & Ba, 2015) using a loss that combines cross-entropy with prompt alignment regularizers (eq. (28)). Full implementation details are provided in Appendix E.3.

Competitors. We compare CNSP against the state-of-the-art prompt-based continual learning methods, including L2P (Wang et al., 2022b), Dual-Prompt (Wang et al., 2022a), CodaPrompt (Smith et al., 2023), CPrompt (Gao et al., 2024), EvoPrompt (Kurniawan et al., 2024), PGP (Qiao et al., 2024), LGCL (Khan et al., 2023), ConvPrompt (Roy et al., 2024), CPG (Lu et al., 2025), and NSP² (Lu et al., 2024).

5.2 MAIN RESULTS AND ANALYSIS

Overall Comparison on ImageNet-21K ViT-B/16. As shown in table 1, CNSP consistently outperforms prior prompt-based continual learning methods on all benchmarks using the ImageNet-21K pretrained ViT-B/16 backbone. On 10-Split DomainNet, for example, CNSP achieves a 9.7% relative reduction in forgetting (7.33% \rightarrow 6.62%) while also improving accuracy from 84.17% to 84.51% over NSP², indicating that enforcing principled null-space consistency yields measurable gains even in accuracy-saturated regimes. Task-wise curves (fig. 3) show the same pattern: CNSP maintains higher or comparable ACC and consistently lower forgetting across all tasks. CNSP also surpasses strong baselines including ConvPrompt and CPG, demonstrating that the proposed matrix-form analysis and right-side projection deliver robust, generalizable benefits in both within-domain and cross-domain continual learning.

Table 1: Comparison with existing prompt-based methods using ImageNet-21K pretrained backbones. Results (%) are reported as mean \pm standard deviation over three runs. Results marked with †, ‡, and § are reproduced by Qiao et al. (2024), Gao et al. (2024), and Lu et al. (2025), respectively, due to the lack of official results.

Method	10-Split CIFAR100		20-Split CIFAR100		10-Split ImageNet-R		10-Split DomainNet	
	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)
L2P	83.83 \pm 0.04	7.63 \pm 0.30	81.29 \pm 0.43 [†]	8.96 \pm 0.38 [†]	61.57 \pm 0.66	9.73 \pm 0.47	81.17 \pm 0.83 [‡]	8.98 \pm 1.25 [‡]
DualPrompt	86.51 \pm 0.33	5.16 \pm 0.09	82.98 \pm 0.47 [†]	8.20 \pm 0.08 [†]	68.13 \pm 0.49	4.68 \pm 0.20	81.70 \pm 0.78 [‡]	8.04 \pm 0.31 [‡]
CodaPrompt	86.25 \pm 0.74	1.67 \pm 0.26	-	-	75.45 \pm 0.56	1.64 \pm 0.10	80.04 \pm 0.79 [‡]	10.16 \pm 0.35 [‡]
CPrompt	87.82 \pm 0.21	5.06 \pm 0.50	83.97 \pm 0.31 [§]	6.85 \pm 0.43 [§]	77.14 \pm 0.11	5.97 \pm 0.68	82.97 \pm 0.34	7.45 \pm 0.93
EvoPrompt	87.97 \pm 0.30	2.60 \pm 0.42	84.64 \pm 0.14	3.98 \pm 0.24	76.83 \pm 0.08	2.78 \pm 0.06	79.50 \pm 0.29 [§]	3.81 \pm 0.36 [§]
PGP	86.92 \pm 0.05	5.35 \pm 0.19	83.74 \pm 0.01	7.91 \pm 0.15	69.34 \pm 0.05	4.53 \pm 0.04	80.41 \pm 0.25 [§]	8.39 \pm 0.18 [§]
LGCL	87.23 \pm 0.21	5.10 \pm 0.15	-	-	69.46 \pm 0.04	4.20 \pm 0.06	-	-
ConvPrompt	88.87 \pm 0.33	4.75 \pm 0.15	87.22 \pm 0.42 [§]	5.43 \pm 0.29 [§]	77.86 \pm 0.25	4.33 \pm 0.24	79.47 \pm 0.35 [§]	6.49 \pm 0.43 [§]
CPG	90.63 \pm 0.44	3.98 \pm 0.65	88.08 \pm 0.77	5.20 \pm 0.64	78.63 \pm 0.52	7.18 \pm 0.62	83.21 \pm 0.67	7.09 \pm 0.82
NSP ²	90.70 \pm 0.17	4.20 \pm 0.35	88.59 \pm 0.46	4.45 \pm 0.60	79.17 \pm 0.63	5.06 \pm 1.07	84.17 \pm 0.54	7.33 \pm 0.94
CNSP (Ours)	91.01 \pm 0.10	3.91 \pm 0.47	89.42 \pm 0.58	4.23 \pm 0.51	79.69 \pm 0.32	4.62 \pm 0.66	84.51 \pm 0.41	6.62 \pm 0.26

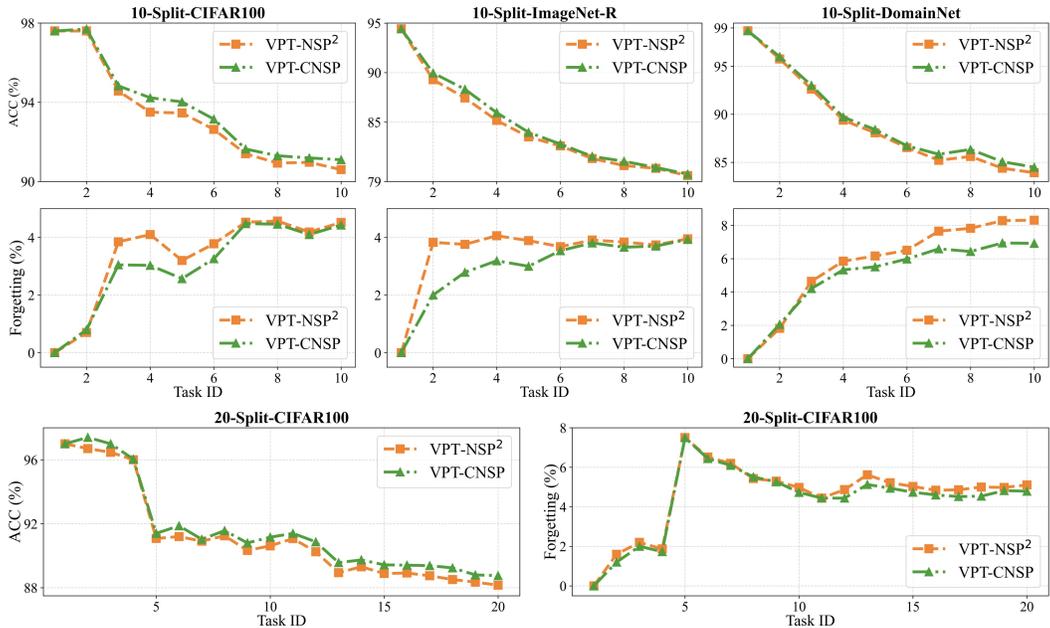


Figure 3: Task-wise performance curves on 10-Split CIFAR100, 20-Split CIFAR100, 10-Split ImageNet-R, and 10-Split DomainNet using the ViT-B/16-21K backbone. For each benchmark, results include the last average accuracy (ACC, %), where higher is better, and the corresponding forgetting (%), where lower indicates better retention of prior knowledge.

Robustness Across Backbones. We further evaluate CNSP on two distinct backbones—ViT-L/16 (AugReg) (Steiner et al., 2021) and the self-supervised DINOv2 ViT-B/14 (Oquab et al., 2023)—with results shown in tables 2 and 3. Across both models, CNSP consistently outperforms NSP². The improvement is especially pronounced on ViT-L/16, where approximation errors in high-dimensional token spaces accumulate, while CNSP’s exact right-side projection remains stable. On 10-Split ImageNet-R, forgetting drops from 6.50% to 4.29% (34% relative reduction), alongside an ACC increase from 83.52% to 84.78%. The gains on DINOv2 further show that CNSP does not rely on supervised ImageNet-21K pretraining and transfers well to self-supervised features. Overall, these results demonstrate that CNSP is robust to variations in model scale and pretraining strategy.

Table 2: Results on four benchmarks with ViT-L/16 (AugReg, IN-21k) backbone.

Method	10-Split ImageNet-R		10-Split DomainNet		10-Split CIFAR100		20-Split CIFAR100	
	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)
Upper Bound	88.58	0.00	90.48	0.00	95.06	0.00	95.06	0.00
Lower Bound	79.92	15.51	80.85	15.88	88.91	10.11	87.40	11.32
NSP ²	83.52	6.50	84.38	10.28	92.20	4.14	90.83	6.34
CNSP (Ours)	84.78	4.29	85.63	8.03	92.45	4.01	91.02	5.24

Table 3: Results on four benchmarks with DINOv2 ViT-B/14 backbone.

Method	10-Split ImageNet-R		10-Split DomainNet		10-Split CIFAR100		20-Split CIFAR100	
	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)
Upper Bound	89.58	0.00	89.65	0.00	93.22	0.00	93.22	0.00
Lower Bound	78.42	17.66	68.87	31.89	78.01	22.27	77.84	21.80
NSP ²	85.22	4.91	84.52	10.06	89.43	5.08	87.15	5.92
CNSP (Ours)	85.67	4.49	85.27	8.50	89.87	4.79	87.29	5.72

Robustness to Domain Shifts. To evaluate robustness under task heterogeneity, we test CNSP on 6-Split DomainNet, where each task corresponds to a distinct visual domain. As shown in table 4, forgetting decreases from 6.67% to 5.84% on run 1 (a 12% relative reduction) with an ACC gain (83.55% \rightarrow 84.32%). These results indicate that CNSP maintains stable performance under substantial domain shifts without relying on assumptions about invariant prompt distributions.

Table 4: Results on the 6-Split DomainNet benchmark across three runs with ViT-B/16-21K.

Method	Run 1		Run 2		Run 3	
	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)	ACC(\uparrow)	Forgetting(\downarrow)
Upper Bound	91.68	0.00	91.94	0.00	91.88	0.00
Lower Bound	63.34	34.11	79.17	18.52	78.22	17.65
NSP ²	83.55	6.67	84.85	7.91	84.76	3.36
CNSP (Ours)	84.32	5.84	85.34	6.44	85.28	3.07

Ablation Study. A detailed ablation study is provided in Appendix E.5. In practice, it shows that the right-side null-space projection is the primary driver of CNSP’s performance, while the variance-preservation loss mainly contributes to training stability.

Summary. CNSP delivers consistent gains across model scales, pretraining paradigms, and domain-shift scenarios. Additional analysis on prompt expressive capacity is provided in Appendix E.6.

6 CONCLUSION

We introduced CNSP, a principled framework that provides the first end-to-end theoretical formulation of prompt-based continual learning. By modeling MHSA, LayerNorm, and classification heads in exact matrix form, CNSP derives unified sufficient consistency conditions for joint feature and head preservation, leading directly to a simple and tractable projection rule. Across diverse benchmarks and Transformer backbones—including large-scale, cross-domain, and self-supervised settings—CNSP consistently improves accuracy and reduces forgetting. Looking forward, CNSP opens several promising directions, including theoretically grounded head designs, extensions to multimodal continual learning, adaptive prompt allocation, and tighter analyses of nonlinear components such as softmax. Additional discussions and limitations are provided in Appendix F.

540 REPRODUCIBILITY STATEMENT

541

542 Our experiments are conducted on publicly available datasets (e.g., CIFAR-100, ImageNet). We
 543 report the details of our training setup, including hyperparameters (batch size, learning rate, opti-
 544 mizer), number of epochs, and evaluation metrics, in Appendix E. To ensure reproducibility, we will
 545 release the complete source code and training/evaluation scripts upon acceptance of the paper.

546

547 REFERENCES

548

549 Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment:
 550 Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF conference on*
 551 *computer vision and pattern recognition*, pp. 113–123, 2019.

552 Danruo Deng, Guangyong Chen, Jianye Hao, Qiong Wang, and Pheng-Ann Heng. Flattening sharp-
 553 ness for dynamic gradient projection memory benefits continual learning. *Advances in Neural*
 554 *Information Processing Systems*, 34:18710–18721, 2021.

555

556 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 557 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
 558 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
 559 *arXiv:2010.11929*, 2020.

560 Zhanxin Gao, Jun Cen, and Xiaobin Chang. Consistent prompting for rehearsal-free continual learn-
 561 ing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 562 pp. 28463–28473, 2024.

563

564 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
 565 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.

566 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
 567 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A criti-
 568 cal analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international*
 569 *conference on computer vision*, pp. 8340–8349, 2021.

570

571 Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.

572 Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and
 573 Ser-Nam Lim. Visual prompt tuning. In *European conference on computer vision*, pp. 709–727.
 574 Springer, 2022.

575

576 Muhammad Gul Zain Ali Khan, Muhammad Ferjad Naeem, Luc Van Gool, Didier Stricker, Fed-
 577 erico Tombari, and Muhammad Zeshan Afzal. Introducing language guidance in prompt-based
 578 continual learning. In *Proceedings of the IEEE/CVF international conference on computer vision*,
 579 pp. 11463–11473, 2023.

580 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
 581 *Conference on Learning Representations*, 2015.

582 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 583 2009.

584

585 Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter
 586 Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with
 587 dirichlet calibration. *Advances in neural information processing systems*, 32, 2019.

588 Muhammad Rifki Kurniawan, Xiang Song, Zhiheng Ma, Yuhang He, Yihong Gong, Yang Qi, and
 589 Xing Wei. Evolving parameterized prompt memory for continual learning. In *Proceedings of the*
 590 *AAAI Conference on Artificial Intelligence*, volume 38, pp. 13301–13309, 2024.

591

592 Guoliang Lin, Hanlu Chu, and Hanjiang Lai. Towards better plasticity-stability trade-off in incre-
 593 mental learning: A simple linear connector. In *Proceedings of the IEEE/CVF conference on*
computer vision and pattern recognition, pp. 89–98, 2022a.

- 594 Sen Lin, Li Yang, Deliang Fan, and Junshan Zhang. Trgp: Trust region gradient projection for
595 continual learning. *arXiv preprint arXiv:2202.02931*, 2022b.
596
- 597 Yue Lu, Shizhou Zhang, De Cheng, Yinghui Xing, Nannan Wang, Peng Wang, and Yanning Zhang.
598 Visual prompt tuning in null space for continual learning. *Advances in neural information pro-
599 cessing systems*, 37:7878–7901, 2024.
- 600 Yue Lu, Shizhou Zhang, De Cheng, Guoqiang Liang, Yinghui Xing, Nannan Wang, and Yanning
601 Zhang. Training consistent mixture-of-experts-based prompt generator for continual learning.
602 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 19152–19160,
603 2025.
- 604 Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The
605 sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165.
606 Elsevier, 1989.
607
- 608 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
609 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
610 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 611 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
612 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
613 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
614
- 615 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching
616 for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference
617 on computer vision*, pp. 1406–1415, 2019.
- 618 Jingyang Qiao, zhizhong zhang, Xin Tan, Chengwei Chen, Yanyun Qu, Yong Peng, and Yuan
619 Xie. Prompt gradient projection for continual learning. In *The Twelfth International Confer-
620 ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=EH2O3h7sBI)
621 [EH2O3h7sBI](https://openreview.net/forum?id=EH2O3h7sBI).
- 622
- 623 Anurag Roy, Riddhiman Moulick, Vinay K Verma, Saptarshi Ghosh, and Abir Das. Convolutional
624 prompting meets language models for continual learning. In *Proceedings of the IEEE/CVF con-
625 ference on computer vision and pattern recognition*, pp. 23616–23626, 2024.
- 626 Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning.
627 *CoRR*, abs/2103.09762, 2021. URL <https://arxiv.org/abs/2103.09762>.
- 628
- 629 James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim,
630 Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual de-
631 composed attention-based prompting for rehearsal-free continual learning. In *Proceedings of the
632 IEEE/CVF conference on computer vision and pattern recognition*, pp. 11909–11919, 2023.
- 633 Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas
634 Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv
635 preprint arXiv:2106.10270*, 2021.
- 636
- 637 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
638 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-
639 tion processing systems*, 30, 2017.
- 640 Cheng Wang. Calibration in deep learning: A survey of the state-of-the-art. *arXiv preprint
641 arXiv:2308.01222*, 2023.
- 642
- 643 Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual
644 learning: Theory, method and application. *IEEE Transactions on Pattern Analysis and Machine
645 Intelligence*, 46(8):5362–5383, 2024.
- 646
- 647 Shipeng Wang, Xiaorong Li, Jian Sun, and Zongben Xu. Training networks in null space of feature
covariance for continual learning. In *Proceedings of the IEEE/CVF conference on Computer
Vision and Pattern Recognition*, pp. 184–193, 2021.

648 Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren,
649 Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for
650 rehearsal-free continual learning. In *European conference on computer vision*, pp. 631–648,
651 2022a.

652 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vin-
653 cent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Pro-
654 ceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 139–149,
655 2022b.

656
657 Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning of context-dependent pro-
658 cessing in neural networks. *Nature Machine Intelligence*, 1(8):364–372, 2019.

659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

APPENDIX

A PROOFS OF LEMMAS AND PROPOSITIONS

Lemma 1. Consider the computation of LayerNorm $\text{LN}(\cdot)$. Let $\mathbf{A} = [\mathbf{A}^{(1)}; \mathbf{A}^{(2)}]$ with $\mathbf{A}^{(1)} \in \mathbb{R}^{n \times d}$ and $\mathbf{A}^{(2)} \in \mathbb{R}^{m \times d}$. Then $\text{LN}(\mathbf{A}) = [\text{LN}(\mathbf{A}^{(1)}); \text{LN}(\mathbf{A}^{(2)})]$.

Proof. Expanding the LayerNorm operations gives

$$\text{LN}(\mathbf{A}) = \frac{\mathbf{A} - \boldsymbol{\mu}_{\mathbf{A}}}{\boldsymbol{\sigma}_{\mathbf{A}}} \cdot \boldsymbol{\gamma} + \boldsymbol{\beta}, \quad \boldsymbol{\mu}_{\mathbf{A}}, \boldsymbol{\sigma}_{\mathbf{A}} \in \mathbb{R}^{n+m}, \quad \boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^d, \quad (31)$$

where broadcasting is applied to match the dimensions $(n+m) \times d$. By matrix reformulation of the broadcast, let

$$\mathbf{C} = \mathbf{I}_d - \frac{1}{d} \mathbf{1}_d \mathbf{1}_d^\top \in \mathbb{R}^{d \times d}, \quad \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}} = \text{diag}(\boldsymbol{\sigma}_{\mathbf{A}}) \in \mathbb{R}^{(n+m) \times (n+m)}, \quad \boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma}) \in \mathbb{R}^{d \times d}. \quad (32)$$

Then

$$\text{LN}(\mathbf{A}) = \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}^{-1} \mathbf{A} \mathbf{C} \boldsymbol{\Gamma} + \mathbf{1}_{n+m} \boldsymbol{\beta}^\top. \quad (33)$$

In practice, the standard deviation of the i -th row is computed as

$$\sqrt{\frac{1}{d} \text{Var}(\mathbf{A}[i]) + \epsilon}, \quad (34)$$

where $\text{Var}(\cdot)$ is the operation to compute variance, $\epsilon > 0$ ensures non-degeneracy. Hence, $\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}$ is invertible. Since $\boldsymbol{\sigma}_{\mathbf{A}}$ corresponds row-wise to \mathbf{A} , we may write

$$\boldsymbol{\sigma}_{\mathbf{A}} = \begin{bmatrix} \boldsymbol{\sigma}_{\mathbf{A}^{(1)}} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\sigma}_{\mathbf{A}^{(2)}} \end{bmatrix}, \quad \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}} = \begin{bmatrix} \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}^{(1)}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}^{(2)}}} \end{bmatrix}. \quad (35)$$

Therefore,

$$\text{LN}(\mathbf{A}) = \begin{bmatrix} \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}^{(1)}}}^{-1} \mathbf{A}^{(1)} \mathbf{C} \boldsymbol{\Gamma} + \mathbf{1}_n \boldsymbol{\beta}^\top \\ \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}^{(2)}}}^{-1} \mathbf{A}^{(2)} \mathbf{C} \boldsymbol{\Gamma} + \mathbf{1}_m \boldsymbol{\beta}^\top \end{bmatrix} = \begin{bmatrix} \text{LN}(\mathbf{A}^{(1)}) \\ \text{LN}(\mathbf{A}^{(2)}) \end{bmatrix}. \quad (36)$$

□

Lemma 2. For $\mathbf{A} \in \mathbb{R}^{n \times m}$ with left null space $\mathcal{N}_L(\mathbf{A}) = \{\mathbf{x} \mathbf{A} = \mathbf{0} \mid \mathbf{x}^\top \in \mathbb{R}^n\}$, it holds that $\mathcal{N}_L(\mathbf{A}) = \mathcal{N}_L(\mathbf{A} \mathbf{A}^\top)$.

Proof. $\forall \mathbf{x} \in \mathcal{N}_L(\mathbf{A})$, we have $\mathbf{x} \mathbf{A} = \mathbf{0}$, which implies $\mathbf{x} \mathbf{A} \mathbf{A}^\top = \mathbf{0}$, i.e., $\mathbf{x} \in \mathcal{N}_L(\mathbf{A} \mathbf{A}^\top)$. Conversely, $\forall \mathbf{x} \in \mathcal{N}_L(\mathbf{A} \mathbf{A}^\top)$, we have $\mathbf{x} \mathbf{A} \mathbf{A}^\top = \mathbf{0}$. Then, $0 = \mathbf{x} \mathbf{A} \mathbf{A}^\top \mathbf{x}^\top = \|\mathbf{A}^\top \mathbf{x}^\top\|^2$, which implies $\mathbf{x} \mathbf{A} = \mathbf{0}$, i.e., $\mathbf{x} \in \mathcal{N}_L(\mathbf{A})$. □

Lemma 3. Consider the operation $\text{MLP}(\cdot)$ in a standard ViT MLP (Dosovitskiy et al., 2020), containing two layers with a GELU non-linearity, defined as a linear projection followed by a GELU activation and another linear projection. For a block matrix $\mathbf{A} = [\mathbf{A}^{(1)}; \mathbf{A}^{(2)}]$, where $\mathbf{A}^{(1)} \in \mathbb{R}^{n \times d}$ and $\mathbf{A}^{(2)} \in \mathbb{R}^{m \times d}$, we have

$$\text{MLP}(\mathbf{A}) = \begin{bmatrix} \text{MLP}(\mathbf{A}^{(1)}) \\ \text{MLP}(\mathbf{A}^{(2)}) \end{bmatrix} \quad (37)$$

Proof. By definition,

$$\text{MLP}(\mathbf{A}) = \phi(\mathbf{A} \mathbf{W}_1 + \mathbf{1}_{n+m} \mathbf{b}_1^\top) \mathbf{W}_2 + \mathbf{1}_{n+m} \mathbf{b}_2^\top, \quad (38)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d \times d_h}$, $\mathbf{W}_2 \in \mathbb{R}^{d_h \times d}$ are weight matrices of two linear projections respectively, $\mathbf{b}_1 \in \mathbb{R}^{d_h}$, $\mathbf{b}_2 \in \mathbb{R}^d$ are bias vectors, and ϕ denotes the element-wise GELU.

Since linear transformations, element-wise nonlinearities, and bias additions act independently on each row, the block structure is preserved. Explicitly,

$$\phi(\mathbf{A} \mathbf{W}_1 + \mathbf{1}_{n+m} \mathbf{b}_1^\top) = \phi\left(\begin{bmatrix} \mathbf{A}^{(1)} \\ \mathbf{A}^{(2)} \end{bmatrix} \mathbf{W}_1 + \begin{bmatrix} \mathbf{1}_n \\ \mathbf{1}_m \end{bmatrix} \boldsymbol{\beta}_1^\top\right) = \begin{bmatrix} \phi(\mathbf{A}^{(1)} \mathbf{W}_1 + \mathbf{1}_n \boldsymbol{\beta}_1^\top) \\ \phi(\mathbf{A}^{(2)} \mathbf{W}_1 + \mathbf{1}_m \boldsymbol{\beta}_1^\top) \end{bmatrix}. \quad (39)$$

Similarly, multiplying by \mathbf{W}_2 and adding bias yields

$$\phi(\mathbf{A}\mathbf{W}_1 + \mathbf{1}_{n+m}\mathbf{b}_1^\top)\mathbf{W}_2 + \mathbf{1}_{n+m}\mathbf{b}_2^\top = \begin{bmatrix} \phi(\mathbf{A}^{(1)}\mathbf{W}_1 + \mathbf{1}_n\mathbf{b}_1^\top)\mathbf{W}_2 + \mathbf{1}_n\mathbf{b}_2^\top \\ \phi(\mathbf{A}^{(2)}\mathbf{W}_1 + \mathbf{1}_m\mathbf{b}_1^\top)\mathbf{W}_2 + \mathbf{1}_m\mathbf{b}_2^\top \end{bmatrix}, \quad (40)$$

which implies that

$$\text{MLP}(\mathbf{A}) = \begin{bmatrix} \text{MLP}(\mathbf{A}^{(1)}) \\ \text{MLP}(\mathbf{A}^{(2)}) \end{bmatrix} \quad (41)$$

□

Proposition 1. Consider two tasks \mathcal{T}_t and \mathcal{T}_{t+1} . If the image-token part of attention outputs remains unchanged for all blocks in VPT, i.e.,

$$\mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(l)}[:N+1] = \mathbf{O}_{\mathbf{Z}_{t,t}}^{(l)}[:N+1], \quad l = 1, 2, \dots, L, \quad (6)$$

then the image-token output of the final block is preserved:

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[:N+1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[:N+1], \quad (7)$$

which we refer to as **feature preservation**. Here $\mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(l)}$, $\mathbf{O}_{\mathbf{Z}_{t,t}}^{(l)}$, $\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}$, and $\mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}$ are computed by eq. (1). In particular, the final [CLS] token is preserved:

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[:1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[:1]. \quad (8)$$

Proof. Without loss of generality, we first omit the task indices and focus on the l -th transformer block. As shown in fig. 1, within a block, $\mathbf{Z}^{(l)}$ is added to the output $\mathbf{O}_{\mathbf{Z}}^{(l)}$ via a residual connection, followed by LayerNorm and an MLP with an additional residual connection:

$$\mathbf{B}_{\mathbf{Z}}^{(l)} = \mathbf{Z}^{(l)} + \mathbf{O}_{\mathbf{Z}}^{(l)}, \quad \mathbf{E}_{\mathbf{Z}}^{(l)} = \mathbf{B}_{\mathbf{Z}}^{(l)} + \text{MLP}\left(\text{LN}\left(\mathbf{B}_{\mathbf{Z}}^{(l)}\right)\right). \quad (42)$$

Since

$$\mathbf{E}_{\mathbf{Z}}^{(l)}[:N+1] = \mathbf{B}_{\mathbf{Z}}^{(l)}[:N+1] + \text{MLP}\left(\text{LN}\left(\mathbf{B}_{\mathbf{Z}}^{(l)}\right)\right)[:N+1], \quad (43)$$

by lemmas 1 and 3, we obtain

$$\text{MLP}\left(\text{LN}\left(\mathbf{B}_{\mathbf{Z}}^{(l)}\right)\right)[:N+1] = \text{MLP}\left(\text{LN}\left(\mathbf{B}_{\mathbf{Z}}^{(l)}[:N+1]\right)\right). \quad (44)$$

Moreover,

$$\mathbf{B}_{\mathbf{Z}}^{(l)}[:N+1] = \mathbf{Z}^{(l)}[:N+1] + \mathbf{O}_{\mathbf{Z}}^{(l)}[:N+1] = \mathbf{X}^{(l)} + \mathbf{O}_{\mathbf{Z}}^{(l)}[:N+1], \quad (45)$$

thus,

$$\mathbf{E}_{\mathbf{Z}}^{(l)}[:N+1] = \mathbf{X}^{(l)} + \mathbf{O}_{\mathbf{Z}}^{(l)}[:N+1] + \text{MLP}\left(\text{LN}\left(\mathbf{X}^{(l)} + \mathbf{O}_{\mathbf{Z}}^{(l)}[:N+1]\right)\right). \quad (46)$$

In VPT, before passing to the next block, the prompt embeddings are discarded and replaced by new learnable prompts:

$$\mathbf{X}^{(l+1)} = \mathbf{E}_{\mathbf{Z}}^{(l)}[:N+1], \quad \mathbf{Z}^{(l+1)} = [\mathbf{X}^{(l+1)}; \mathbf{P}^{(l+1)}]. \quad (47)$$

Now consider two tasks \mathcal{T}_t and \mathcal{T}_{t+1} . We prove the statement by induction on l :

(i) When $l = 1$, by assumption and setting, we have

$$\mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(1)}[:N+1] = \mathbf{O}_{\mathbf{Z}_{t,t}}^{(1)}[:N+1], \quad \mathbf{X}_{t,t+1}^{(1)} = \mathbf{X}_{t,t}^{(1)} = \mathbf{X}_t, \quad (48)$$

where \mathbf{X}_t is the input image tokens of VPT from task \mathcal{T}_t dataset. According to eq. (46), it follows that

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(1)}[:N+1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(1)}[:N+1], \quad \mathbf{X}_{t,t+1}^{(2)} = \mathbf{X}_{t,t}^{(2)}. \quad (49)$$

(ii) Suppose that for some $l = k (< L)$,

$$\mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(k)}[:N+1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(k)}[:N+1], \quad (50)$$

810 which implies

$$811 \mathbf{X}_{t,t+1}^{(k+1)} = \mathbf{X}_{t,t}^{(k+1)}. \quad (51)$$

812 For $l = k + 1$, we know

$$813 \mathbf{O}_{\mathbf{Z}_{t,t+1}}^{(k+1)}[: N + 1] = \mathbf{O}_{\mathbf{Z}_{t,t}}^{(k+1)}[: N + 1]. \quad (52)$$

814 By eqs. (46) and (51), it follows that

$$815 \mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(k+1)}[: N + 1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(k+1)}[: N + 1]. \quad (53)$$

816 (iii) By induction, for all $l = 1, 2, \dots, L$, we have

$$817 \mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(l)}[: N + 1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(l)}[: N + 1]. \quad (54)$$

818 Especially, the final [CLS] token fed into the classification head is unchanged, i.e.,

$$819 \mathbf{E}_{\mathbf{Z}_{t,t+1}}^{(L)}[: 1] = \mathbf{E}_{\mathbf{Z}_{t,t}}^{(L)}[: 1]. \quad (55)$$

820 □

821 B MATRIX-FORM CHARACTERIZATION OF LAYERNORM

822 LayerNorm operates row-wise on tokens and normalizes each embedding using its own mean and
823 standard deviation. While this definition is straightforward in implementation, the usual broad-
824 casting notation obscures the algebraic structure of the operation and complicates the analysis of
825 consistency conditions. We therefore present an exact matrix-form characterization that is exactly
826 equivalent to the standard implementation and preserves its semantics.

827 Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a sequence of token embeddings. Standard LayerNorm is defined as

$$828 \text{LN}(\mathbf{A}) = \frac{\mathbf{A} - \boldsymbol{\mu}_{\mathbf{A}}}{\boldsymbol{\sigma}_{\mathbf{A}}} \cdot \boldsymbol{\gamma} + \boldsymbol{\beta}, \quad \boldsymbol{\mu}_{\mathbf{A}}, \boldsymbol{\sigma}_{\mathbf{A}} \in \mathbb{R}^n, \quad \boldsymbol{\gamma}, \boldsymbol{\beta} \in \mathbb{R}^d, \quad (56)$$

829 where $\boldsymbol{\mu}_{\mathbf{A}}, \boldsymbol{\sigma}_{\mathbf{A}}$ are row-wise statistics and $\boldsymbol{\gamma}, \boldsymbol{\beta}$ are the learned affine parameters applied element-
830 wise along the feature dimension.

831 **Exact representation of mean subtraction.** Introduce the centering matrix

$$832 \mathbf{C} = \mathbf{I}_d - \frac{1}{d} \mathbf{1}_d \mathbf{1}_d^\top \in \mathbb{R}^{d \times d}, \quad (57)$$

833 which yields

$$834 (\mathbf{A}\mathbf{C})[i, j] = \mathbf{A}[i, j] - \left(\frac{1}{d} \mathbf{A} \mathbf{1}_d \mathbf{1}_d^\top\right)[i, j] = \mathbf{A}[i, j] - \frac{1}{d} \sum_{k=1}^d \mathbf{A}[i, k] = \mathbf{A}[i, j] - \boldsymbol{\mu}_{\mathbf{A}}[i]. \quad (58)$$

835 Thus, the mean subtraction step is absorbed into the linear operator \mathbf{C} and no information is lost.

836 **Exact representation of variance scaling.** Define the diagonal matrix

$$837 \mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}} = \text{diag}(\boldsymbol{\sigma}_{\mathbf{A}}) \in \mathbb{R}^{n \times n}. \quad (59)$$

838 Left multiplication by $\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}^{-1}$ implements the row-wise division:

$$839 (\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}^{-1} \mathbf{A})[i, j] = \frac{\mathbf{A}[i, j]}{\boldsymbol{\sigma}_{\mathbf{A}}[i]}. \quad (60)$$

840 (The addition of a small positive constant in LayerNorm ensures invertibility of $\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}$.) Therefore,
841 the standard deviation division is exactly realized by the linear operator $\mathbf{D}_{\boldsymbol{\sigma}_{\mathbf{A}}}^{-1}$.

842 **Affine transformation in matrix form.** Define

$$843 \boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma}) \in \mathbb{R}^{d \times d}, \quad \mathbf{1}_n \boldsymbol{\beta}^\top \in \mathbb{R}^{n \times d}. \quad (61)$$

844 Then for every i, j :

$$845 (\mathbf{A}\boldsymbol{\Gamma})[i, j] = \sum_{k=1}^d \mathbf{A}[i, k] \boldsymbol{\Gamma}[k, j] = \mathbf{A}[i, j] \boldsymbol{\gamma}[j], \quad (62)$$

864 showing that right-multiplication with Γ is exactly equivalent to the broadcasted element-wise affine
865 transformation.

866 **Complete matrix-form LayerNorm.** The LayerNorm operation can thus be written in exact matrix
867 form:

$$868 \text{LN}(A) = D_{\sigma_A}^{-1} A C \Gamma + \mathbf{1}_n \beta^\top. \quad (63)$$

869 This representation strictly preserves the semantics of standard LayerNorm while revealing its lin-
870 ear-affine structure and removing the ambiguity of implicit broadcasting, and it enables the consis-
871 tency constraints in Sec. 4 to be derived in closed algebraic form.

872 C CONSISTENT NULL-SPACE PROJECTION IN PRACTICE

873 In the main text (section 4.3), we derived sufficient conditions for feature preservation, requiring that
874 each prompt update ΔP lies in the left null space of the constraint matrix R . Theoretically, this is
875 enforced by projecting ΔP onto the null space of RR^\top (see eq. (30)). However, in practice, due to
876 finite-precision arithmetic, exact zero singular values rarely occur, making it nontrivial to determine
877 the effective nullity. To address this issue, following Lu et al. (2024), we estimate the nullity of
878 RR^\top using an *adaptive thresholding criterion* based on the inflection point of the singular value
879 spectrum. Specifically, we compute the discrete second derivative of sorted singular values and use
880 its maximizer to determine the cut-off between the null and non-null components. Formally, let λ_j
881 denote the j -th singular value of RR^\top (sorted in descending order). The adaptive estimate of the
882 nullity is given by

$$883 \text{Nullity}(RR^\top) = D - \arg \max_j \{\lambda_{j+1} - 2\lambda_j + \lambda_{j-1}\}_{j=2}^{D-1}, \quad (64)$$

884 where the criterion detects the point of maximum curvature change in the spectrum. Empirically,
885 this approach has proven more robust than fixed thresholds in VPT Lu et al. (2024), providing a
886 numerically stable approximation of the null space dimension.

887 With the estimated null space, we then construct an orthogonal basis U_0 for $\mathcal{N}_L(RR^\top)$ and form
888 the projection matrix $B_R = U_0 U_0^\top$. During training, the raw prompt update ΔP_{raw} from back-
889 propagation is consistently projected onto this null space, ensuring that the sufficient conditions for
890 feature preservation are satisfied at every optimization step. This procedure constitutes our *Consis-*
891 *istent Null-Space Projection (CNSP)* mechanism.

892 The detailed training algorithm is summarized in Algorithm 1, and the computation of the null-space
893 projection matrix is given in Algorithm 2.

894 D COMPARISON AND DISCUSSION WITH NSP²

895 CNSP and NSP² both rely on projection-based updates for prompt tuning in Transformers, yet they
896 differ substantially in how thoroughly they characterize and enforce representational consistency.
897 Below, we highlight the key distinctions in a concise and analytically oriented manner.

898 **Multi-head Derivation.** NSP² first derives sufficient conditions under a single-head setting and
899 then extends them to multi-head attention by direct concatenation, which lacks a full derivation and
900 leaves potential inter-head interactions unexamined. CNSP provides a complete multi-head analysis
901 by explicitly unfolding the MHSA forward computation. This yields rigorous per-head sufficient
902 conditions and clarifies the role of the shared output projection W_o . The resulting formulation
903 precisely characterizes when head-wise invariance guarantees MHSA-level representational preser-
904 vation, offering a fully grounded multi-head theory.

905 **LayerNorm Modeling and Distributional Assumptions.** NSP² treats LayerNorm’s centering and
906 scaling operations as cancelable scalar factors, simplifying its algebraic structure at the cost of omit-
907 ting key broadcasted operations. CNSP introduces an explicit matrix-form characterization of Lay-
908 erNorm, including:

- 909 • the centering matrix C ,
- 910 • the per-sample normalization matrix D_σ^{-1} ,

Algorithm 1 CNSP for VPT in Continual Learning

Inputs: Datasets $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{|\mathcal{T}_t|}$ for task $\mathcal{T}_t \in \{\mathcal{T}_1, \mathcal{T}_2, \dots\}$; ViT model $f_{\text{ViT}}(\cdot | \{\mathbf{P}_t^{(l)}\}_{l=1}^L)$ with the prompts $\{\mathbf{P}_t^{(l)}\}_{l=1}^L$ to be optimized (the updates of classification head is omitted for simplicity); the gram matrix of constraints \mathbf{G} ; the projection matrix \mathbf{B}_R

Outputs: Optimized prompts $\mathbf{P}_t^{(l)}$

- 1: **Initialization:** Randomly initialize $\mathbf{P}_t^{(l)}$; set $\mathbf{G}^{(l)} = \mathbf{0}$, $\mathbf{B}_R^{(l)} = \mathbf{I}$
- 2: **for** task $\mathcal{T}_t \in \{\mathcal{T}_1, \mathcal{T}_2, \dots\}$ **do**
- 3: **repeat**
- 4: Sample a mini-batch $\mathbf{x}_t, y_t \sim \mathcal{D}_t$
- 5: $\hat{y}_t \leftarrow f_{\text{ViT}}(\mathbf{x}_t | \{\mathbf{P}_t^{(l)}\}_{l=1}^L)$
- 6: $\mathcal{L}_{\text{total}} \leftarrow \text{CrossEntropy}(\hat{y}_t, y_t)$ ▷ the classification loss
- 7: **if** $t > 1$ **then**
- 8: Compute \mathcal{L}_{std} by eq. (28) ▷ loss of prompts standard deviation
- 9: $\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{total}} + \mathcal{L}_{\text{std}}$
- 10: **end if**
- 11: Get raw prompts update $\Delta \mathbf{P}_{\text{raw}}^{(l)}$ from optimizer using $\mathcal{L}_{\text{total}}$
- 12: **if** $t > 1$ **then**
- 13: $\Delta \mathbf{P}^{(l)} \leftarrow \Delta \mathbf{P}_{\text{raw}}^{(l)} \mathbf{B}_R^{(l)}$ ▷ consistent null-space projection, eq. (30)
- 14: **else**
- 15: $\Delta \mathbf{P}^{(l)} \leftarrow \mathbf{P}_{\text{raw}}^{(l)}$
- 16: **end if**
- 17: $\mathbf{P}_t^{(l)} \leftarrow \mathbf{P}_t^{(l)} - lr \times \Delta \mathbf{P}^{(l)}$
- 18: **until** convergence
- 19: **if** $t = 1$ **then**
- 20: $\mathbf{R}_v^{(l)} \leftarrow \mathbf{C} \Gamma^{(l)} \mathbf{W}_v^{(l)}$ ▷ cache the second constraint matrix
- 21: **end if**
- 22: **for** $\mathbf{x}_t^{(i)} \in \mathcal{D}_t$ **do**
- 23: $\mathbf{R}_k^{(l)} \leftarrow \mathbf{C} \Gamma^{(l)} \mathbf{W}_k^{(l)} \mathbf{Q}_{\mathbf{x}_t^{(i)}}^{(l)\top}$ ▷ by forward propagation $f_{\text{ViT}}(\mathbf{x}_t^{(i)} | \{\mathbf{P}_t^{(l)}\}_{l=1}^L)$
- 24: $\mathbf{G}^{(l)} \leftarrow \mathbf{G}^{(l)} + \mathbf{R}_k^{(l)} \mathbf{R}_k^{(l)\top} + \mathbf{R}_v^{(l)} \mathbf{R}_v^{(l)\top}$
- 25: **end for**
- 26: $\mathbf{B}_R^{(l)} \leftarrow \text{ComputingNullSpaceProjection}(\mathbf{G}^{(l)})$ ▷ using algorithm 2
- 27: **end for**

Algorithm 2 Computing Null-Space Projection Matrix**Inputs:** Gram matrix \mathbf{G} **Outputs:** Null-space projection matrix \mathbf{B}_R

- 1: $\mathbf{U}, \Sigma, _ \leftarrow \text{SVD}(\mathbf{G})$ ▷ sorted by the singular values in descending order
- 2: Computing the nullity N_G of \mathbf{G} by eq. (64)
- 3: $\mathbf{U}_0 \leftarrow \mathbf{U}[D - N_G : D]$ ▷ get the left singular vectors of zero singular values
- 4: $\mathbf{B}_R \leftarrow \frac{\mathbf{U}_0 \mathbf{U}_0^\top}{\|\mathbf{U}_0 \mathbf{U}_0^\top\|_F}$ ▷ improve numerical stability by its Frobenius norm

- the affine scaling matrix Γ .

This exact formulation captures the full behavior of LayerNorm and enables a variance-only consistency condition, relaxing the stronger mean–variance assumption employed in NSP². The relaxation improves stability while maintaining theoretical correctness.

Implementation of Sufficient Conditions. A key practical distinction lies in how the two methods translate theoretical constraints into implementable projections.

NSP² requires two SVDs per layer: (i) a left projection: $M \times M$ SVD (where M is the number of prompt tokens), and (ii) a right projection: $D \times D$ SVD (where D is the feature dimension). CNSP’s

972 unified derivation leads to a single right-side $D \times D$ projection, significantly reducing computational
 973 cost and memory overhead.

974 Because CNSP models LayerNorm exactly, the resulting constraints naturally admit a right-sided
 975 nullification form that avoids circular dependencies on statistics such as $D_\sigma^{(-1)}$. This makes the pro-
 976 jection stable across tasks and smoothly scalable to larger backbones (e.g., ViT-L/16), as observed
 977 in our experiments.
 978

979 **Classification Head Preservation.** NSP² analyzes invariance within intermediate layers but does
 980 not address how classification heads interact with updated prompts. CNSP extends representational
 981 consistency to the end-to-end pipeline by deriving sufficient conditions for preserving both the back-
 982 bone features (MHSA + LayerNorm) and the classification head. This provides the first mathemat-
 983 ical characterization of when task-specific heads remain valid after prompt updates and how they
 984 decompose from feature preservation constraints.

985 **Summary.** Overall, NSP² provides a valuable and influential starting point for principled prompt-
 986 based continual learning. CNSP addresses the theoretical gaps in NSP² by supplying rigorous multi-
 987 head derivations, an exact LayerNorm formulation, relaxed distributional assumptions, and the first
 988 end-to-end consistency analysis. Its unified projection rule results in a more stable and scalable
 989 implementation that preserves task performance while maintaining theoretical soundness.
 990

991 E EXPERIMENTS

992 E.1 DATASETS AND BENCHMARKS

993 To thoroughly evaluate the proposed CNSP method in prompt-based continual learning scenarios,
 994 we conduct experiments on five continual learning benchmarks spanning both class-incremental and
 995 domain-incremental settings:
 996
 997

- 998 • **10/20-Split CIFAR-100** (Krizhevsky et al., 2009): CIFAR-100 consists of 60,000 32×32
 999 color images across 100 classes, with 500 training images and 100 test images per class.
 1000 Following common practice, it is randomly partitioned into 10 or 20 tasks, where each task
 1001 contains 10 or 5 classes, respectively.
- 1002 • **10-Split ImageNet-R** (Hendrycks et al., 2021): ImageNet-R contains 30,000 images from
 1003 200 ImageNet classes, covering diverse styles such as art, cartoons, and sketches. It is
 1004 randomly divided into 10 tasks, each consisting of 20 classes.
- 1005 • **10-Split DomainNet** (Peng et al., 2019): DomainNet is a cross-domain dataset covering
 1006 345 everyday object classes across six domains (clipart, real, sketch, infograph, painting,
 1007 and quickdraw). Since the number of images per class varies significantly, we follow the
 1008 existing setting (Gao et al., 2024; Lu et al., 2024) by selecting the 200 largest classes and
 1009 randomly splitting them into 10 tasks, each containing 20 classes with samples drawn from
 1010 multiple domains.
- 1011 • **6-Split DomainNet** (Peng et al., 2019): To evaluate robustness under strong distribution
 1012 shifts, we construct a domain-incremental protocol where each task corresponds to one of
 1013 the six visual domains in DomainNet (Real, Sketch, Painting, Clipart, Infograph, Quick-
 1014 draw). For each domain, we select the 20 classes with the largest sample counts to form
 1015 representative tasks.
 1016

1017 E.2 METRICS

1018 In continual learning, evaluation is typically based on the accuracy matrix $\mathbf{A} \in \mathbb{R}^{T \times T}$, where T is
 1019 the total number of tasks and the entry $A_{i,j}$ denotes the accuracy on task j after training on task i .
 1020 We report two standard metrics: Last Average Accuracy and Last Average Forgetting.
 1021

- 1022 • **Last Average Accuracy (ACC)** measures the overall performance on all tasks after com-
 1023 pleting the training sequence:

$$1024 \text{ACC} = \frac{1}{T} \sum_{j=1}^T A_{T,j},$$

- **Last Average Forgetting (Forgetting)** quantifies the average performance drop on past tasks after learning subsequent ones:

$$\text{Forgetting} = \frac{1}{T-1} \sum_{j=1}^{T-1} \max_{i \in \{1, \dots, T-1\}} (A_{i,j} - A_{T,j}).$$

E.3 IMPLEMENTATION DETAILS

We follow the class-incremental learning protocol, where task classes are disjoint and the task identity is unknown at inference. Our main experiments are conducted with ViT-B/16 (Dosovitskiy et al., 2020) pretrained on ImageNet-21K as the backbone. In addition, we also report results on two other backbones: ViT-L/16 pretrained with AugReg (Steiner et al., 2021) and the self-supervised DINOv2 ViT-B/14 (Oquab et al., 2023). Each transformer layer is augmented with four learnable prompts. The prompt embeddings are shared across tasks and initialized from a uniform distribution in $[-1, 1]$.

Models are trained with Adam (Kingma & Ba, 2015) (learning rate of 0.01, a batch size of 256, $\beta_1 = 0.9$, $\beta_2 = 0.999$, weight decay 5×10^{-5}). Each task is trained for 10 epochs on all benchmarks, which is sufficient for convergence (as shown in Appendi E.4). NSP² is reproduced from its official codebase and hyperparameters, with the sole modification of training for 10 epochs per task instead of the originally reported 100. Empirical evidence in Appendix E.4 confirms fully convergence, disentangling forgetting from underfitting. The training objective combines cross-entropy loss with a feature standard-deviation alignment loss for prompt embeddings (see eq. (28)). The alignment weight is set to $\lambda = 1.0$ for 10-Split CIFAR-100, 20-Split CIFAR-100, 10-Split ImageNet-R, 10-Split DomainNet, and to $\lambda = 2.0$ for 6-Split DomainNet. The temperature parameter in cross-entropy is tuned via cross-validation and set to 28, 25, 30, 30, and 30 for 10-Split CIFAR-100, 20-Split CIFAR-100, 10-Split ImageNet-R, 10-Split DomainNet, and 6-Split DomainNet, respectively.

E.4 CONVERGENCE GUARANTEE

To justify the training protocol adopted in our experiments, we report the training accuracy and loss curves of both NSP² and CNSP across tasks on four benchmarks (10-Split-CIFAR100, 20-Split-CIFAR100, 10-Split-DomainNet, and 10-Split-ImageNet-R). As shown in figs. 5 to 8, all curves exhibit rapid convergence within the first 5 epochs: training accuracy quickly increases and stabilizes, while training loss consistently decreases and plateaus. By epoch 10, both methods have fully converged across tasks and datasets, as further confirmed by the shaded regions indicating the last four epochs and the black curve representing the average accuracy/loss across all tasks at each epoch. These results provide empirical evidence that training each task for 10 epochs is sufficient to ensure convergence, without sacrificing performance or stability.

E.5 ABLATION STUDY

Table 5 reports both the overall comparison and the ablation of CNSP. The lower bound (naïve VPT without continual learning mechanisms) suffers from severe forgetting, while joint training serves as an oracle upper bound. NSP² provides strong performance, yet CNSP consistently achieves higher ACC and lower forgetting across all four benchmarks, substantially narrowing the gap toward the oracle while clearly surpassing naïve tuning.

The ablation study further clarifies each component’s role. Removing the variance preservation loss (\mathcal{L}_{std}) causes moderate ACC drops and increased forgetting, suggesting that it contributes mainly to stability. In contrast, removing the right-side projection (\mathbf{B}_R) leads to pronounced degradation on all benchmarks, confirming it as the key mechanism for effective feature preservation. Finally, the VPT-only baseline performs worst overall, highlighting the substantial gains achieved by CNSP beyond a frozen backbone.

E.6 ANALYSIS OF PROMPT EXPRESSIVE CAPACITY

To assess whether the null-space projection restricts the expressive capacity of prompts as the number of tasks increases, we analyze the evolution of the prompt-update subspace over tasks. Specifi-

Table 5: Overall comparison and ablation study of CNSP on four benchmarks by removing the variance preservation loss (\mathcal{L}_{std}), the right-side null-space projection (\mathbf{B}_R), and the entire CNSP module (“Lower Bound”). Upper Bound indicates training with access to all task data simultaneously.

Method	10-Split CIFAR100		20-Split CIFAR100		10-Split ImageNet-R		10-Split DomainNet	
	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)	ACC(↑)	Forgetting(↓)
Upper Bound	93.59	0.00	93.59	0.00	84.62	0.00	89.52	0.00
NSP ²	90.60	4.52	88.12	5.11	79.60	3.95	83.92	8.31
CNSP(Ours)	91.10	4.43	88.76	4.79	79.75	3.93	84.49	6.92
w/o \mathcal{L}_{std}	90.58	5.17	88.22	6.13	79.72	4.30	84.29	7.58
w/o \mathbf{B}_R	85.01	14.42	84.80	14.84	73.98	17.62	78.29	19.85
Lower Bound	83.72	16.17	83.88	15.55	73.33	18.33	74.30	24.94

cally, we examine the rank of the projection matrices $\{\mathbf{B}_t^{(l)}\}_{l=1}^L$, where $\mathbf{B}_t^{(l)}$ denotes the projection matrix used for task t in layer l . These matrices determine the allowable update directions at each layer l . We define the average rank ratio across layers for each incremental task t (> 1) as:

$$\rho_t = \frac{1}{L} \sum_{l=1}^L \frac{\text{rank}(\mathbf{B}_t^{(l)})}{D},$$

where D is the token embedding dimension. A high value of ρ_t indicates that most update directions remain available, implying that the expressive capacity of prompts is preserved.

Table 6 reports the average rank ratio across tasks for three benchmarks using ViT-B/16-21K. Across all setting, ρ_t remains extremely high ($>97.52\%$) and shows no degradation as tasks progress. These results demonstrate that CNSP does not collapse the expressive capacity of prompts. Instead, the null-space constraint removes only the directions that would violate feature preservation, while maintaining a near-full-rank update subspace and strong adaptability throughout continual learning.

Table 6: Average rank ratio (%) of projection matrices across layers for each task.

Benchmark	ρ_2	ρ_3	ρ_4	ρ_5	ρ_6	ρ_7	ρ_8	ρ_9	ρ_{10}
10-Split ImageNet-R	97.53	97.53	97.53	97.54	97.53	97.53	97.54	97.53	97.53
10-Split CIFAR100	97.53	97.54	97.54	97.54	97.54	97.54	97.54	97.55	97.55
10-Split DomainNet	97.52	97.54	97.53	97.53	97.53	97.54	97.53	97.53	97.53

F DISCUSSION AND LIMITATIONS

F.1 DISCUSSION ON HEAD PRESERVATION

In CNSP, head preservation (eq. (9)) naturally ensured under the standard protocol of using task-specific classification heads in prompt-based continual learning. As illustrated in fig. 4, during training only the head associated with the current task is updated, while all previously learned heads remain frozen. Combined with the feature-preservation condition established in our framework, this implies that every previously trained head receives invariant representations and therefore retains its original input-output mapping, directly satisfying proposition 2.

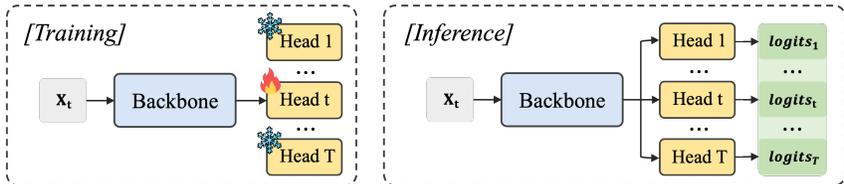
Beyond the theoretical formulation, the practical cost of employing multiple task-specific heads is negligible. In a ViT-B/16 ($\approx 86\text{M}$ parameters), prompts (four per layer) contribute $\approx 0.037\text{M}$ and a 20-way classifier $\approx 0.015\text{M}$. Even with ten tasks, the cumulative head parameters remain below 0.2% of the backbone, indicating that the multi-head design does not materially increase model size.

1134 Importantly, CNSP itself does not rely on any particular choice of classifier head. Because our
 1135 consistency formulation cleanly separates feature preservation from head behavior, the framework is
 1136 compatible with a wide range of head architectures, including prototype-based, cosine-normalized,
 1137 low-rank, or other more compact or shared designs. In this work, we adopt the standard independent
 1138 linear heads primarily to align with common practice in prompt-based continual learning and to
 1139 ensure a fair comparison with existing baselines, rather than due to any methodological constraint
 1140 imposed by CNSP.

1141 At inference, we adopt the common practice of concatenating logits from all task-specific heads and
 1142 selecting the class with the highest confidence. This procedure can be interpreted probabilistically:
 1143 each head produces calibrated class scores over its own label subset, and concatenation yields a
 1144 unified score vector over the global label space, upon which a global decision is made via arg max.

1145 A practical concern is that logits produced by different heads may not be calibrated to the same
 1146 scale. While this issue is independent of the CNSP consistency theory, it may affect fairness across
 1147 tasks in the concatenation step. To mitigate this, we apply temperature scaling (Guo et al., 2017)
 1148 during training to smooth predictions and improve cross-head comparability. More advanced cali-
 1149 bration techniques (e.g., Dirichlet calibration (Kull et al., 2019)) or normalization-based heads (e.g.,
 1150 cosine/prototype classifiers) can be readily integrated into CNSP without modifying the underlying
 1151 consistency formulation.

1152 Finally, when tasks contain visually similar classes, multiple heads may assign high confidence to
 1153 the same example. Addressing such conflicts—e.g., via logit normalization, lightweight gating, or
 1154 shared prototype heads—presents an interesting avenue for future work. Importantly, these strategies
 1155 operate on top of CNSP’s theory rather than altering it, enabled by the unified feature and head
 1156 preservation guarantees established in our framework.



1164 Figure 4: Training and inference of task-specific heads. During training, only the current task head
 1165 is updated while others remain frozen. During inference, all heads run in parallel and their logits are
 1166 concatenated, with the final prediction given by the most confident class, eliminating the need for
 1167 task identity.
 1168

1170 F.2 DISCUSSION ON PROMPT POSITION

1171 In practice, we adopt the strategy of appending prompts after the image patches, following NSP² (Lu
 1172 et al., 2024), to ensure direct comparability. From lemmas 1 and 3, it follows that the preservation
 1173 conditions in propositions 1 and 2 are invariant to the relative position of prompt tokens within the
 1174 input sequence. In other words, whether prompts are prepended before the [CLS] token or appended
 1175 after the patch embeddings, the algebraic form of the preservation constraints remains unchanged.
 1176 This invariance arises from the row-wise independence of LayerNorm and MLP operations, which
 1177 treat each token identically regardless of ordering. While these properties hold in theory, a sys-
 1178 tematic empirical investigation into how different prompt insertion strategies influence performance
 1179 preservation is beyond the scope of this work. We consider this an interesting and important direc-
 1180 tion for future research.

1182 F.3 FUTURE DIRECTIONS AND LIMITATIONS

1183 While CNSP provides a unified theoretical view of feature and head preservation, several aspects
 1184 remain outside the scope of this work. First, although our analysis gives exact treatments of Lay-
 1185 erNorm and MHSA, we keep the softmax-induced constraints analytically tractable by not fully re-
 1186 laxing all nonlinear interactions; consequently, our sufficient conditions remain conservative rather
 1187 than approaching necessary-and-sufficient characterizations. Second, we adopt task-specific classi-

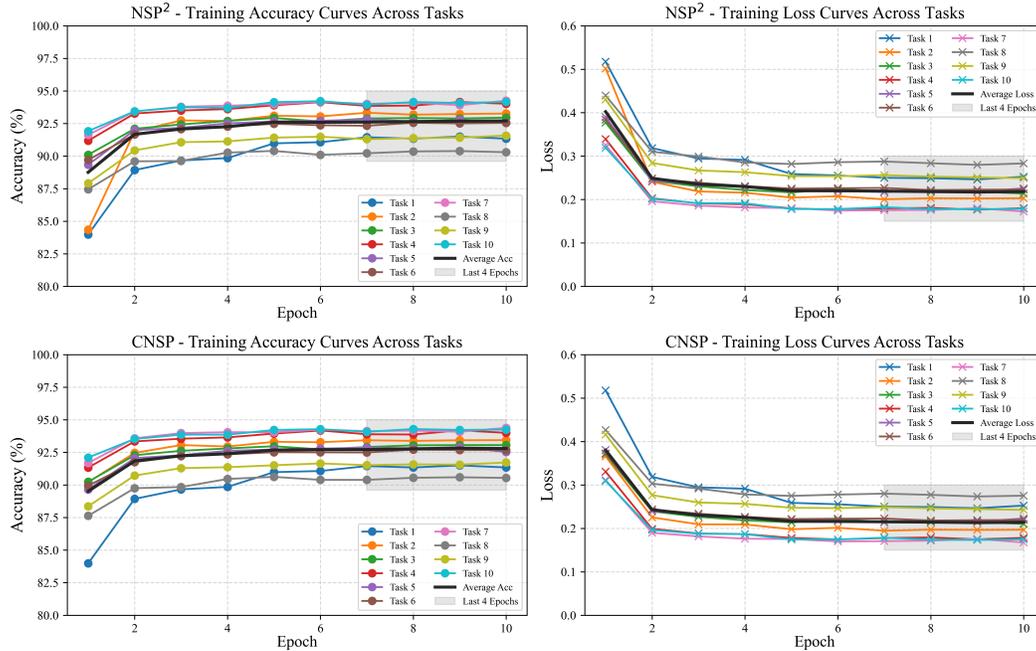
1188 fication heads for theoretical clarity; scaling to highly heterogeneous tasks may benefit from shared
 1189 or normalized head designs. Third, our experiments focus on single-modality vision benchmarks.
 1190

1191 These scope boundaries point toward several promising extensions. Natural next steps include
 1192 multi-modal continual learning with vision–language models, adaptive prompt allocation under task
 1193 heterogeneity, and tighter theoretical analysis of softmax and other nonlinear components within
 1194 Transformer-based continual learning.
 1195

1196 G LLM USAGE STATEMENT

1197 We acknowledge the use of Large Language Models (LLMs) as an assistive tool in the preparation
 1198 of this work. Their role was limited to grammar refinement, LaTeX formatting assistance, language
 1199 polishing, and minor code debugging/modification. No part of the research methodology, experi-
 1200 ments, or data analysis was delegated to LLMs. All conceptual contributions, design of experiments,
 1201 and interpretation of results were carried out by the authors. The authors take full responsibility for
 1202 all content presented in this paper, including any text or code that may have been refined with the
 1203 assistance of LLMs.
 1204
 1205

1206 **10-Split-CIFAR100**



1228 Figure 5: Training accuracy and loss curves across 10 tasks on the 10-Split-CIFAR100 benchmark.
 1229 The top row reports results of NSP², and the bottom row reports CNSP. Left: training accuracy
 1230 across tasks. Right: training loss across tasks. Each colored curve corresponds to one task, while
 1231 the black curve denotes the average across all tasks at each epoch.
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

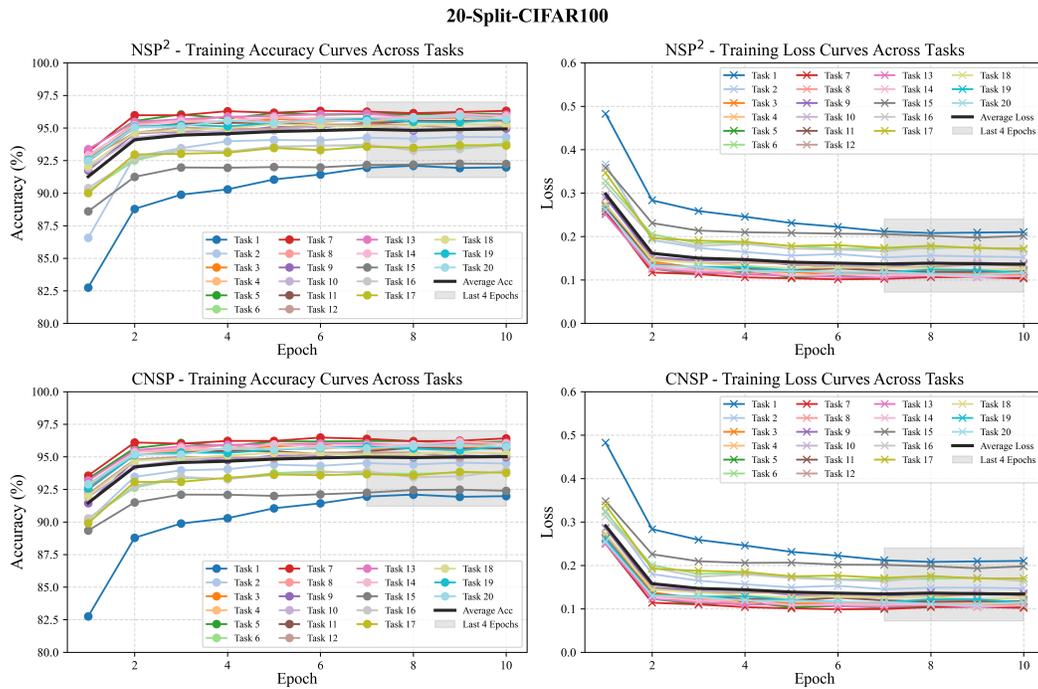


Figure 6: Training accuracy and loss curves across 20 tasks on the 20-Split-CIFAR100 benchmark. The top row reports results of NSP², and the bottom row reports CNSP. Left: training accuracy across tasks. Right: training loss across tasks. Each colored curve corresponds to one task, while the black curve denotes the average across all tasks at each epoch.

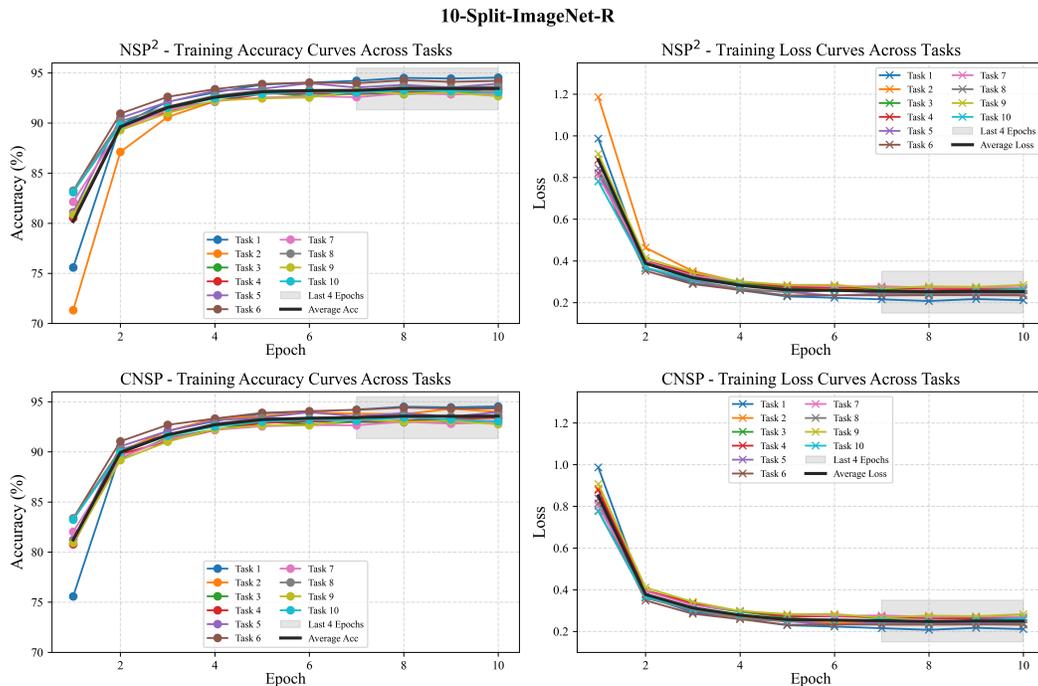


Figure 7: Training accuracy and loss curves across 10 tasks on the 10-Split-ImageNet-R benchmark. The top row reports results of NSP², and the bottom row reports CNSP. Left: training accuracy across tasks. Right: training loss across tasks. Each colored curve corresponds to one task, while the black curve denotes the average across all tasks at each epoch.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

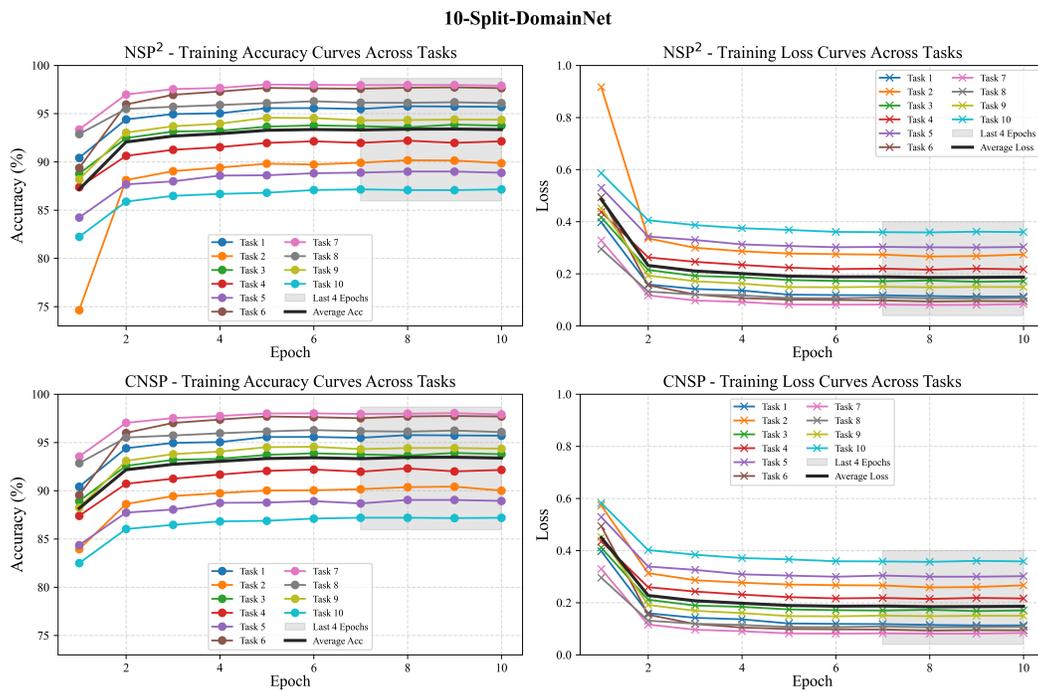


Figure 8: Training accuracy and loss curves across 10 tasks on the 10-Split-DomainNet benchmark. The top row reports results of NSP², and the bottom row reports CNSP. Left: training accuracy across tasks. Right: training loss across tasks. Each colored curve corresponds to one task, while the black curve denotes the average across all tasks at each epoch.