

# Supplementary Material

## 1 Computational Overhead

Module/Eq.	Complexity/ $O(\cdot)$	Total FLOPs	Peak RAM
Graph Construction (1)	$N_c^2 \cdot D$	$7.41 \times 10^{14}$	<b>778 MB</b>
Spectral Embedding (2)	$N_c^3$	$2.10 \times 10^{12}$	163 KB
Cost Matrix (3)	$M_c \cdot N_c \cdot d$	$6.30 \times 10^{12}$	655 KB
Hungarian Match	$M_c \cdot N_c^2$	$8.40 \times 10^{13}$	741 KB
PD-CFD Loss (9)	$(N_c + M_c)N_{freq}d$	$9.20 \times 10^{12}$	1.03 MB
<b>Total FAST</b>	–	$8.43 \times 10^{14}$	$\approx 782$ MB
ResNet-50 Comparison	–	$5.25 \times 10^{15}$	$\approx 6.25$ GB

Table 1: Expensive Computation Breakdown on ImageNet-1K.

**Note:** FAST assumes per-class selection. ResNet-50 FLOPs are for one epoch inference on ImageNet-1K with a batch size of 256. Total feature extraction overhead of FAST on ImageNet-1K is 16% of ResNet-50.

Keep Rate	10%	Time(h)	20%	Time(h)	30%	Time(h)	Full
previous SOTA <sup>‡</sup>	53.6	7.9	62.6	21.1	64.5	30.2	71.2
FAST <sup>‡</sup>	<b>56.2</b>	<b>2.5</b>	<b>63.9</b>	<b>6.1</b>	<b>66.9</b>	<b>8.2</b>	71.2

Table 2: Performance Comparison on ImageNet-1K. <sup>‡</sup>On CPU.

## 2 Optimization Stability

The *Hungarian mapping* evolves from dynamic exploration to stable anchoring, ensuring robust continuous-to-discrete alignment. Despite discrete updates to  $\pi$ , the loss maintains continuity at cost-equalizing *Voronoi boundaries*, and divergence is prevented by topological constraints that limit gradient fluctuations to the local manifold neighborhood. Since *Hungarian algorithm* minimizes the transport cost globally, gradient descent on the smooth objective  $\tilde{Y}$  ensures  $\mathcal{L}(\tilde{Y}^{(t+1)}, \pi^{(t+1)}) \leq \mathcal{L}(\tilde{Y}^{(t)}, \pi^{(t)})$ , yielding monotonic descent optimization:  $0 \leq \mathcal{L}_{total}^{(t+1)} \leq \mathcal{L}_{total}^{(t)}$ , where *Monotone Convergence Theorem* guarantees convergence.

## 3 Necessity of Embeddings

In LLM experiments, we utilize general Sentence-BERT embeddings as standard preprocessing for manifold construction. The structural distinction between the encoder-only embedder and decoder-only LLM introduces a general semantic prior rather than coupled architectural bias, confirmed by Table ?? that shows FAST’s robust generalization.

Model	L-FAST	L-Rand	Q-FAST	Q-Rand	M-FAST	M-Rand
Acc (%)	39.0 <sup>↑8.7</sup>	35.9 <sup>↑5.6</sup>	72.3 <sup>↑6.1</sup>	70.9 <sup>↑4.7</sup>	63.1 <sup>↑4.0</sup>	61.6 <sup>↑2.5</sup>

Table 3: Supplementary Experiments on LLM datasets.

**Notation:** L: LLaMA2-7B, Q: Qwen2.5-7B, M: Mistral-7B, Rand:Random. <sup>↑</sup> indicates the relative improvement to the base model.

## 4 More Ablation Studies

Use CIFAR-10 dataset (10% KR, training ResNet-18) unless noted; random sampling shows  $75.70\%_{\pm 7.63}$  accuracy.

Method	Acc (%)	Time (s)	Energy (Wh)	Std	Opt/Steps
Pixel-Opt <sup>‡</sup>	61.52	3770	29.75	13.73	1300
ResNet-50 <sup>†</sup>	81.51	627.67	37.75	3.32	230
FAST <sup>‡</sup>	<b>90.32</b>	<b>353</b>	<b>1.409</b>	<b>1.21</b>	<b>80</b>

Table 4: Graph Feature Extractor Ablation on <sup>‡</sup>CPU, <sup>†</sup>GPU.

**Analysis:** Pixel-level optimization fails due to the sparsity and noise of high-dimensional space, while incurring massive computational costs (10× slower). ResNet-50 features suffer from architectural bias and computational overhead.

Parameter	Ours	$k = 5$	$k = 50$	$d = 8$	$d = 128$
Trustworthiness	<b>0.95<math>\pm 0.01</math></b>	0.73 $\pm 0.10$	0.52 $\pm 0.09$	0.44 $\pm 0.09$	0.92 $\pm 0.02$
Continuity	0.91 $\pm 0.03$	0.49 $\pm 0.15$	0.93 $\pm 0.06$	0.92 $\pm 0.03$	0.87 $\pm 0.06$
Acc (%)	<b>90.32<math>\pm 1.21</math></b>	72.12 $\pm 5.56$	85.01 $\pm 3.27$	78.86 $\pm 3.13$	86.73 $\pm 4.16^*$

Table 5: Hyperparameter Ablation (neighbor scale  $k$  and reduced dimension  $d$ ). \*Time doubles. FAST adopts  $k = 15$  and  $d = 32$ .

**Analysis:** Small  $k$  fractures the manifold, while large  $k$  causes collapse, connecting distant classes; low  $d$  bottlenecks information, and high  $d$  adds noise and doubles optimization time. Our settings optimally preserve topology.

Strategy	DAS	PES	PUS	US	Collinear	PDAS
Acc (%)	88.89 $\pm 2.23$	87.12 $\pm 2.60$	88.66 $\pm 2.01$	86.60 $\pm 2.73$	74.16 $\pm 7.77$	<b>90.32<math>\pm 1.21</math></b>
Opt/Steps	150	400	250	430	100*	<b>80</b>

Table 6: Impact of Frequency Selection. \*Optimization failure mode. Refer to paper Fig. 9 for illustrations.

**Analysis:** Baselines unstably converge due to suboptimal frequency selection; PDAS achieves optimal stability by progressively selecting discriminative frequencies.

Method	CIFAR-10	RESISC-45	DTD
FAST-Vanilla	88.17 $\pm 2.12$	82.01 $\pm 3.99$	41.15 $\pm 3.20$
FAST-PD	<b>90.32<math>\pm 1.21</math></b>	<b>85.00<math>\pm 2.01</math></b>	<b>45.77<math>\pm 2.53</math></b>

Table 7: Phase-Decoupled CFD Ablation (full ablation in Fig. 7).

**Analysis:** *CFD* fails to capture high-frequency details due to amplitude-phase coupling. *PD-CFD* resolves this, showing substantial gains on texture-rich DTD (+4.62%) and RESISC-45 (+2.99%) compared to CIFAR-10 (+2.15%).

## 5 Additional Analysis on Limitations of KL and CE for High-Order Moment Alignment

**Notation.** Let  $P$  be an unknown target distribution on  $\mathcal{X} \subseteq \mathbb{R}^d$  with density  $p$ . An exponential family is

$$\mathcal{Q}_T := \left\{ q_\theta(x) = h(x) \exp(\theta^\top T(x) - A(\theta)) \mid \theta \in \Theta \subset \mathbb{R}^m \right\},$$

where  $T : \mathcal{X} \rightarrow \mathbb{R}^m$  is a vector of sufficient statistics,  $A(\theta) := \log \int h(x) \exp(\theta^\top T(x)) dx$  is the log-partition function (cumulant generating function), and  $h$  is a base density. We write  $\mu(\theta) := \nabla A(\theta) = \mathbb{E}_{q_\theta}[T(X)]$  for the mean-parameter map.

For  $k \in \mathbb{N}$ , define the (truncated) polynomial statistics

$$T^{(\leq k)}(x) = (x, xx^\top, x^{\otimes 3}, \dots, x^{\otimes k}),$$

so that  $\mathcal{Q}_{\leq k}$  denotes the exponential family whose sufficient statistics are all monomials up to degree  $k$  (with a suitable choice of  $h$  to ensure normalizability).

**Proposition 1** (I-projection yields moment matching on the chosen statistics). *Consider the KL minimization (I-projection) of  $P$  onto an exponential family  $\mathcal{Q}_T$ :*

$$\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(P \| Q_\theta) = \arg \min_{\theta} \int p(x) \log \frac{p(x)}{q_\theta(x)} dx. \quad (1)$$

*Then  $A$  is convex and the KL objective is strictly convex in  $\theta$  on the mean-parameter interior; the unique minimizer  $\theta^*$  (when it exists) satisfies the first-order optimality condition*

$$\nabla_{\theta} \left( \mathbb{E}_P[\theta^\top T(X)] - A(\theta) \right) \Big|_{\theta=\theta^*} = \mathbb{E}_P[T(X)] - \nabla A(\theta^*) = 0.$$

*Equivalently,*

$$\mathbb{E}_P[T(X)] = \mathbb{E}_{Q_{\theta^*}}[T(X)].$$

*Proof (via convex duality and Fenchel–Legendre conjugacy).* Expanding the KL divergence,

$$D_{\text{KL}}(P \| Q_\theta) = -H(P) - \mathbb{E}_P[\theta^\top T(X) - A(\theta) + \log h(X)],$$

where  $H(P)$  and  $\mathbb{E}_P[\log h(X)]$  are  $\theta$ -independent constants. Thus minimizing  $D_{\text{KL}}(P \| Q_\theta)$  is equivalent to maximizing the concave functional

$$\mathcal{L}(\theta) := \mathbb{E}_P[\theta^\top T(X)] - A(\theta).$$

Since  $A$  is convex (indeed,  $A$  is the log-moment generating function of  $T$  under  $h$ ),  $\mathcal{L}$  is concave and (under standard interiority conditions) has a unique maximizer  $\theta^*$ . The Fenchel–Legendre conjugate of  $A$  is  $A^*(\mu) := \sup_{\theta} \{\langle \theta, \mu \rangle - A(\theta)\}$ ; strong duality yields that the optimal  $\theta^*$  satisfies  $\nabla A(\theta^*) = \mu^* = \mathbb{E}_P[T(X)]$  (i.e. the moment-matching equations). Finally  $\nabla A(\theta) = \mathbb{E}_{Q_\theta}[T(X)]$  by standard exponential-family calculus, completing the proof.  $\square$

**Corollary 1** (Distributions controlled by KL under restricted sufficient statistics). *Let  $\mathcal{Q}_T$  be any exponential family with sufficient statistics  $T$ . Then the I-projection  $Q_{\theta^*}$  aligns exactly those coordinates of  $T$ :*

$$\mathbb{E}_P[T_i(X)] = \mathbb{E}_{Q_{\theta^*}}[T_i(X)] \quad \text{for each component } T_i,$$

*but imposes no necessary constraint on expectations of functions  $f$  lying outside the linear span of  $\{1, T_1, \dots, T_m\}$ . In particular, if  $T = T^{(\leq k)}$  contains all monomials up to degree  $k$ , then KL minimization guarantees matching of all raw moments up to order  $k$ , and does not in general constrain any  $(k+1)$ -st or higher-order moments.*

*Proof.* Immediate from Proposition 1, noting that the optimality system is exactly the linear system matching  $\mathbb{E}_P[T]$  and  $\mathbb{E}_{Q_\theta}[T]$ . Any  $f$  outside the closed linear span of  $\{1, T\}$  cannot be represented as  $\alpha_0 + \alpha^\top T$ , hence its expectation is unconstrained by the KKT system.  $\square$

**Corollary 2** (Gaussian family ( $k=2$ ) controls mean and covariance but not higher cumulants). *Let  $\mathcal{Q}_{\mathcal{N}} = \{\mathcal{N}(\mu, \Sigma)\}$ . Then  $\mathcal{Q}_{\mathcal{N}}$  is an exponential family with  $T(x) = (x, xx^\top)$  and  $\nabla A(\mu, \Sigma) = (\mathbb{E}[X], \mathbb{E}[XX^\top])$ . The I-projection  $Q_{\theta^*} = \mathcal{N}(\mu^*, \Sigma^*)$  thus satisfies*

$$\mathbb{E}_{Q_{\theta^*}}[X] = \mathbb{E}_P[X], \quad \mathbb{E}_{Q_{\theta^*}}[XX^\top] = \mathbb{E}_P[XX^\top],$$

*but neither skewness (third cumulant) nor kurtosis (fourth cumulant) is constrained or minimized by the KL objective in general.*

*Proof.* A direct specialization of Corollary 1 with  $k=2$ .  $\square$

**Proposition 2** (Augmenting statistics raises the matched moment order, but only up to that order). Let  $\mathcal{Q}_{\leq k}$  denote the exponential family with  $T^{(\leq k)}$ . Then the sequence of I-projections  $\{Q^{(k)}\}_{k \geq 1}$  defined by

$$Q^{(k)} \in \arg \min_{Q \in \mathcal{Q}_{\leq k}} D_{\text{KL}}(P \| Q)$$

satisfies, for each fixed  $k$ ,

$$\mathbb{E}_{Q^{(k)}}[X^{\otimes r}] = \mathbb{E}_P[X^{\otimes r}] \quad \text{for all } r = 1, 2, \dots, k,$$

and there is no general guarantee that  $\mathbb{E}_{Q^{(k)}}[X^{\otimes r}]$  aligns with  $\mathbb{E}_P[X^{\otimes r}]$  for any  $r > k$ .

*Proof.* By Proposition 1 applied to  $T^{(\leq k)}$ , the optimality conditions enforce equality of the first  $k$  raw-moment tensors. Since the KL objective reduces to a linear functional of  $T$  minus  $A(\theta)$ , any statistics not included in  $T$  do not appear in the optimality system and remain uncontrolled.  $\square$

**Remark 1** (Consequences for KL/CE-based evaluation of distributional alignment). Since cross-entropy minimization is equivalent to minimizing  $D_{\text{KL}}(P \| Q_\theta)$  when  $P$  is fixed, Propositions 1–2 imply a fundamental limitation: Under a restricted exponential family, KL/CE can certify alignment of at most the moments encoded in  $T$  (e.g., mean and covariance for Gaussians), and it provides no guarantees for higher-order moments or cumulants that are not present in  $T$ . Hence, unless higher-order statistics are explicitly included in the model (i.e., increasing  $k$ ), KL/CE-based fitting and evaluation do not ensure moment alignment.

**Example 1** (Gaussian vs. heavy-tailed target: kurtosis mismatch at the KL optimum). Let  $P$  be a zero-mean heavy-tailed distribution (e.g., a centered Laplace in 1D) with variance  $\sigma^2$  and kurtosis  $\kappa_P > 3$ . Its I-projection onto  $\mathcal{Q}_{\mathcal{N}}$  is  $Q^* = \mathcal{N}(0, \sigma^2)$ , which matches mean and variance by Corollary 2. However,  $\kappa_{Q^*} = 3 \neq \kappa_P$ , i.e., fourth-order moments are not aligned at the KL optimum—a concrete manifestation of Remark 1.

**Remark 2** (Mixtures vs. single exponential families). A Gaussian mixture model (GMM) is a convex combination of Gaussians and not a single exponential family with a fixed finite-dimensional  $T$ . Therefore, the structural limitation in Propositions 1–2 does not apply in the same form to GMMs: with sufficiently many components, a GMM can approximate higher-order structures arbitrarily well. However, practical optimization remains nonconvex and capacity-limited, so alignment may still fail in practice despite the absence of a finite-dimensional  $T$ .

## A General Counterexample: Small KL yet Divergent Higher-Order Moments

We now give a short, self-contained construction showing that (beyond model restrictions) the KL divergence itself does not control unbounded test-function expectations, such as moments.

**Proposition 3** (KL can be arbitrarily small while  $k$ -th moments diverge). Fix any integer  $k \geq 3$ . There exists a sequence of distributions  $P_M$  and a reference  $Q$  such that

$$D_{\text{KL}}(P_M \| Q) \rightarrow 0 \quad \text{but} \quad |\mathbb{E}_{P_M}[X^k] - \mathbb{E}_Q[X^k]| \rightarrow \infty.$$

*Proof.* Let  $Q = \mathcal{N}(0, 1)$  and  $R_M = \mathcal{N}(M, 1)$ . Define the mixture

$$P_M = (1 - \varepsilon_M)Q + \varepsilon_M R_M, \quad \varepsilon_M = \frac{1}{M^2 \log M}.$$

(i) *KL upper bound by joint convexity.* By convexity of  $D_{\text{KL}}(\cdot \| Q)$  in its first argument,

$$D_{\text{KL}}(P_M \| Q) \leq (1 - \varepsilon_M) \cdot D_{\text{KL}}(Q \| Q) + \varepsilon_M \cdot D_{\text{KL}}(R_M \| Q) = \varepsilon_M \cdot \frac{M^2}{2} = \frac{1}{2 \log M} \rightarrow 0, \quad (2)$$

where  $D_{\text{KL}}(\mathcal{N}(M, 1) \| \mathcal{N}(0, 1)) = \frac{M^2}{2}$ .

(ii) *Divergence of the  $k$ -th moment gap.* For  $k \geq 3$ ,  $\mathbb{E}_{R_M}[X^k] \sim M^k$ . Hence

$$|\mathbb{E}_{P_M}[X^k] - \mathbb{E}_Q[X^k]| = \varepsilon_M |\mathbb{E}_{R_M}[X^k] - \mathbb{E}_Q[X^k]| \gtrsim \varepsilon_M M^k = \frac{M^{k-2}}{\log M} \xrightarrow{M \rightarrow \infty} \infty. \quad (3)$$

Thus, KL can be made arbitrarily small while the  $k$ -th moment discrepancy diverges.  $\square$

**Interpretation.** Pinsker’s inequality controls only *bounded* test functions via total variation; polynomial moments are unbounded, so small KL does not imply closeness of higher moments. Proposition 3 complements the structural limitation (Propositions 1–2) by showing a *metric* limitation: even without any modeling restrictions, KL does not bound higher-order moments.

## 6 Limitations of MMD for Aligning Higher-Order Statistics

Maximum Mean Discrepancy (MMD) is widely used as a nonparametric distributional metric, defined for a positive definite kernel  $k$  by

$$\text{MMD}_k(P, Q) = \|\mu_P - \mu_Q\|_{\mathcal{H}_k} = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} (\mathbb{E}_P[f(X)] - \mathbb{E}_Q[f(X)]), \quad (4)$$

where  $\mathcal{H}_k$  is the reproducing kernel Hilbert space (RKHS) associated with  $k$ . When  $k$  is characteristic, MMD metrizes weak convergence. However, for the purpose of coreset selection—where the objective is to align *all* statistical moments (including heavy-tailed or high-order cumulants)—weak convergence is insufficient. We show below that with the commonly used bounded kernels (e.g., Gaussian, Laplacian), MMD does *not* control higher-order moments, even when  $\text{MMD}_k(P, Q)$  is arbitrarily small.

**Setup.** Let  $k(x, y)$  be any bounded positive definite kernel with  $\sup_{x, y} |k(x, y)| \leq K < \infty$  (this includes the Gaussian RBF). Let  $Q$  be any reference distribution with finite moments of all orders.

We construct the following sequence:

$$P_M = (1 - \varepsilon_M)Q + \varepsilon_M R_M,$$

where  $R_M$  is a distribution concentrated at radius  $\|x\| \approx M$ , and  $\varepsilon_M$  is a vanishing mixing weight.

The construction is analogous to the KL example in Proposition 3 but adapted to the MMD geometry.

**Proposition 4** (Small MMD does not control higher-order moments). *For any bounded kernel  $k$  and any integer  $r \geq 3$ , there exists a sequence  $P_M$  such that*

$$\text{MMD}_k(P_M, Q) \rightarrow 0 \quad \text{but} \quad |\mathbb{E}_{P_M}[X^r] - \mathbb{E}_Q[X^r]| \rightarrow \infty.$$

*Proof.* Because  $k$  is bounded, the RKHS norm of the kernel mean embedding satisfies

$$\|\mu_{P_M} - \mu_Q\|_{\mathcal{H}_k} \leq \mathbb{E}_{P_M, Q}[|k(X, Y)|] \leq K.$$

More precisely,

$$\mu_{P_M} = (1 - \varepsilon_M)\mu_Q + \varepsilon_M \mu_{R_M},$$

hence

$$\text{MMD}_k(P_M, Q) = \varepsilon_M \|\mu_{R_M} - \mu_Q\|_{\mathcal{H}_k} \leq 2K \varepsilon_M.$$

Taking  $\varepsilon_M = 1/\log M$  yields

$$\text{MMD}_k(P_M, Q) \leq \frac{2K}{\log M} \xrightarrow{M \rightarrow \infty} 0.$$

Next, choose  $R_M$  such that almost all its mass lies on  $\|x\| \approx M$ . Then the  $r$ -th moment satisfies

$$\mathbb{E}_{R_M}[X^r] \asymp M^r,$$

and therefore

$$\mathbb{E}_{P_M}[X^r] = (1 - \varepsilon_M)\mathbb{E}_Q[X^r] + \varepsilon_M \mathbb{E}_{R_M}[X^r] \asymp \mathbb{E}_Q[X^r] + \varepsilon_M M^r. \quad (5)$$

Since  $\varepsilon_M M^r = M^r / \log M \rightarrow \infty$  for any  $r \geq 3$ , we obtain

$$|\mathbb{E}_{P_M}[X^r] - \mathbb{E}_Q[X^r]| \rightarrow \infty,$$

even while  $\text{MMD}_k(P_M, Q) \rightarrow 0$ . □

**Interpretation.** This result parallels Proposition 3 for KL: both reveal a *metric limitation* independent of modeling assumptions. For bounded kernels, MMD metrizes weak convergence but cannot control expectations of unbounded test functions—particularly polynomial functions that capture higher-order moments, tail behavior, or heavy-tailed anisotropy.

Consequently, when MMD is used as an objective for coresets construction, moment alignment depends critically on the expressiveness of the chosen kernel. If the kernel is insufficient to probe heavy tails or high-order interactions, the resulting coreset may perfectly match the MMD score yet deviate drastically in higher-order statistics—precisely the regime where our PD-CFD metric provides a more faithful discrepancy measure.

## 7 Characteristic Functions: Complete Statistical Representation and Fourier–Analytic Foundations

**Notation.** Let  $P$  be a Borel probability measure on  $\mathbb{R}^d$  with random vector  $X \sim P$ . The *characteristic function* (CF) of  $P$  is

$$\varphi_P(\omega) := \int_{\mathbb{R}^d} e^{i\omega^\top x} P(dx) = \mathbb{E}_P[e^{i\omega^\top X}], \quad \omega \in \mathbb{R}^d.$$

Each  $\varphi_P$  is bounded ( $|\varphi_P| \leq 1$ ), uniformly continuous, and positive definite.

**Proposition 5** (Fourier inversion: recovering  $P$  from  $\varphi_P$ ). *Suppose  $P$  admits a density  $p \in L^1(\mathbb{R}^d)$ . Then  $p$  can be recovered from its characteristic function by the inverse Fourier transform:*

$$p(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{-i\omega^\top x} \varphi_P(\omega) d\omega, \quad \text{for a.e. } x \in \mathbb{R}^d.$$

More generally, for any  $f \in L^1(\mathbb{R}^d)$  with Fourier transform  $\widehat{f}(\omega) := \int e^{i\omega^\top x} f(x) dx$ ,

$$\int_{\mathbb{R}^d} f(x) P(dx) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \widehat{f}(\omega) \varphi_P(\omega) d\omega.$$

Thus, the knowledge of  $\varphi_P$  on  $\mathbb{R}^d$  uniquely determines all integrals  $\int f dP$ , and hence the measure  $P$  itself.

*Proof (Fourier transform duality).* The Fourier transform  $\mathcal{F} : f \mapsto \widehat{f}$  is a linear bijection between  $L^1$  and bounded continuous functions, satisfying the Plancherel identity  $\int |f|^2 = (2\pi)^{-d} \int |\widehat{f}|^2$  for  $f \in L^2$ . Since  $e^{i\omega^\top x}$  is the kernel of this transform,

$$\int f(x) P(dx) = \int f(x) \left( \frac{1}{(2\pi)^d} \int e^{-i\omega^\top x} \varphi_P(\omega) d\omega \right) dx = \frac{1}{(2\pi)^d} \int \widehat{f}(\omega) \varphi_P(\omega) d\omega. \quad (6)$$

If  $\varphi_P$  is known everywhere, the right-hand side gives  $\int f dP$  for all  $f \in L^1 \cap L^2$ , implying uniqueness of  $P$ .  $\square$

**Proposition 6** (Lévy’s continuity theorem: weak convergence and uniqueness). *Let  $\{P_n\}$  be a sequence of probability measures on  $\mathbb{R}^d$  with characteristic functions  $\varphi_{P_n}$ . Suppose that  $\varphi_{P_n}(\omega) \rightarrow \varphi(\omega)$  for each  $\omega$ , and that the limit  $\varphi$  is itself a characteristic function (i.e. positive definite and  $\varphi(0) = 1$ ). Then  $P_n \Rightarrow P$ , where  $P$  is the distribution with CF  $\varphi$ . In particular, if  $\varphi_P(\omega) \equiv \varphi_Q(\omega)$  for all  $\omega$ , then  $P = Q$ .*

*Proof (weak convergence via Fourier test functions).* For any  $f \in C_c^\infty(\mathbb{R}^d)$ , its Fourier transform  $\widehat{f}$  is rapidly decaying. By Fubini and dominated convergence,

$$\int f(x) P_n(dx) = \frac{1}{(2\pi)^d} \int \widehat{f}(\omega) \varphi_{P_n}(\omega) d\omega \xrightarrow{n \rightarrow \infty} \frac{1}{(2\pi)^d} \int \widehat{f}(\omega) \varphi(\omega) d\omega = \int f(x) P(dx). \quad (7)$$

Hence  $\int f dP_n \rightarrow \int f dP$  for all  $f \in C_c^\infty$ , which implies  $P_n \Rightarrow P$  by the Portmanteau theorem. Taking  $P_n = P$  and  $\varphi_P = \varphi_Q$  gives the uniqueness  $P = Q$ .  $\square$

**Proposition 7** (Bochner’s theorem: positive-definite functions as Fourier transforms of measures). *A continuous function  $\psi : \mathbb{R}^d \rightarrow \mathbb{C}$  with  $\psi(0) = 1$  is the characteristic function of some probability measure on  $\mathbb{R}^d$  if and only if it is positive definite, i.e. for all  $n \in \mathbb{N}$ , all  $\omega_1, \dots, \omega_n \in \mathbb{R}^d$ , and all  $c_1, \dots, c_n \in \mathbb{C}$ ,*

$$\sum_{i,j=1}^n c_i \bar{c}_j \psi(\omega_i - \omega_j) \geq 0.$$

Moreover, for every finite nonnegative measure  $\mu$  on  $\mathbb{R}^d$ , its Fourier transform  $\psi(\omega) := \int e^{i\omega^\top x} d\mu(x)$  is continuous and positive definite. Thus, continuous positive-definite functions with  $\psi(0) = 1$  are in one-to-one correspondence with characteristic functions of probability measures.

*Proof (spectral representation).* For any finite positive measure  $\mu$ ,  $\psi(\omega) = \int e^{i\omega^\top x} d\mu(x)$  is continuous and satisfies the positive-definite inequality above by direct calculation. Conversely, if  $\psi$  is continuous and positive definite, by the classical Bochner–Khinchine theorem there exists a unique finite nonnegative measure  $\mu$  such that  $\psi$  is its Fourier transform. Setting  $\mu(\mathbb{R}^d) = 1$  gives a probability measure  $P = \mu$ , with  $\psi = \varphi_P$ .  $\square$

**Proposition 8** (Matching characteristic functions on  $\mathbb{R}^d$  equals to matching the joint distribution). *Let  $P, Q$  be Borel probability measures on  $\mathbb{R}^d$  with characteristic functions  $\varphi_P, \varphi_Q$ . Define a nonnegative weighting function  $w \in L^1(\mathbb{R}^d)$  satisfying  $\text{supp}(w) = \mathbb{R}^d$  and the frequency-domain distance*

$$\mathcal{D}_w^2(P, Q) := \int_{\mathbb{R}^d} w(\omega) |\varphi_P(\omega) - \varphi_Q(\omega)|^2 d\omega.$$

Then the following are equivalent:

$$\mathcal{D}_w(P, Q) = 0 \iff \varphi_P(\omega) = \varphi_Q(\omega) \forall \omega \iff P = Q. \quad (8)$$

*Proof (Lévy–Bochner synthesis).* If  $\mathcal{D}_w(P, Q) = 0$ , then  $\varphi_P = \varphi_Q$  almost everywhere on  $\{w > 0\}$ . Since both  $\varphi_P, \varphi_Q$  are continuous, equality extends to all  $\omega \in \mathbb{R}^d$ . By Proposition 6,  $\varphi_P \equiv \varphi_Q$  implies  $P = Q$ . Conversely, if  $P = Q$ , then clearly  $\mathcal{D}_w(P, Q) = 0$ . Thus, equality of CFs on the full frequency domain is equivalent to equality of the underlying distributions.  $\square$

Intuitively,  $D_w(P, Q)$  measures the squared discrepancy between characteristic functions across all frequencies, weighted by  $w$ . When  $w$  has full support, this distance captures the complete distributional difference. We now connect this frequency-domain view to kernel methods via Bochner’s theorem.

**Proposition 9** (Kernel formulation via Bochner’s theorem). *Let  $k(x, y) = \kappa(x - y)$  be a bounded, continuous, translation-invariant kernel with spectral measure  $\Lambda$  satisfying*

$$\kappa(t) = \int_{\mathbb{R}^d} e^{i\omega^\top t} d\Lambda(\omega).$$

Then the maximum mean discrepancy (MMD) between  $P, Q$  in the RKHS of  $k$  admits the spectral representation

$$\text{MMD}_k^2(P, Q) = \int_{\mathbb{R}^d} |\varphi_P(\omega) - \varphi_Q(\omega)|^2 d\Lambda(\omega).$$

Moreover, if the kernel is characteristic—equivalently, the support of  $\Lambda$  is all of  $\mathbb{R}^d$ —then  $\text{MMD}_k(P, Q) = 0$  if and only if  $P = Q$ .

*Proof (spectral integration).* Expanding expectations and applying Bochner’s representation,

$$\begin{aligned} \mathbb{E}k(X, Y) &= \int_{\mathbb{R}^d} \mathbb{E}e^{i\omega^\top(X-Y)} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} \varphi_P(\omega) \overline{\varphi_Q(\omega)} d\Lambda(\omega). \end{aligned} \quad (9)$$

Substituting into the definition  $\text{MMD}_k^2(P, Q) = \mathbb{E}k(X, X') + \mathbb{E}k(Y, Y') - 2\mathbb{E}k(X, Y)$  yields the stated formula. If  $\text{MMD}_k(P, Q) = 0$  while  $\Lambda$  has full support, the integrand’s continuity forces  $\varphi_P \equiv \varphi_Q$ , which by Proposition 6 implies  $P = Q$ .  $\square$

**Summary and Connection to CFD.** Propositions 2–6 collectively establish that a probability distribution is uniquely and completely determined by its characteristic function, and that matching characteristic functions over the entire frequency domain—either directly via a weighted  $L^2$  distance  $D_w$  or indirectly via a characteristic kernel MMD—is equivalent to matching the full joint distribution. The Fourier inversion theorem connects CFs to densities, Bochner’s theorem connects positive-definite kernels to spectral measures, and Lévy’s continuity theorem ensures that equality of CFs implies equality of distributions. Consequently, frequency-domain metrics such as our CFD provide a principled way to capture all moments and dependencies, beyond what marginal or covariance-based criteria can express.

## 8 Moments and Cumulants via Taylor Expansion of Characteristic Functions

### 8.1 Preliminaries and Notations

Let  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$  be a random vector with law  $P$ . Its *characteristic function* (CF) is

$$\varphi(\omega) = \mathbb{E}\left[e^{i\omega^\top X}\right], \quad \omega = (\omega_1, \dots, \omega_d) \in \mathbb{R}^d,$$

and the *log-characteristic function* (log-CF) is  $\psi(\omega) = \log \varphi(\omega)$  whenever  $\varphi(\omega) \neq 0$ .

**Multi-index notation.** A *multi-index* is  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}_0^d$  with  $|\alpha| := \sum_{j=1}^d \alpha_j$ ,  $\alpha! := \prod_{j=1}^d \alpha_j!$ ,  $\omega^\alpha := \prod_{j=1}^d \omega_j^{\alpha_j}$ , and  $X^\alpha := \prod_{j=1}^d X_j^{\alpha_j}$ . The mixed partial derivative is

$$\partial^\alpha := \frac{\partial^{|\alpha|}}{\partial \omega_1^{\alpha_1} \dots \partial \omega_d^{\alpha_d}}.$$

**Standing condition (finite moment of order  $|\alpha|$ ).** Throughout, whenever we speak of  $\partial^\alpha \varphi(0)$  we assume  $\mathbb{E}|X^\alpha| < \infty$ . This ensures, via dominated convergence, that differentiation and expectation can be interchanged.

### 8.2 Limitations of Marginal-Only Matching

Classical matching of *univariate* moments (per-coordinate means/variances/skewness/kurtosis) aligns only the marginals  $\{P_{X_j}\}_{j=1}^d$ . However, cross-variable dependence (pairwise, triple-wise, and beyond) lives in *mixed* moments such as  $\mathbb{E}[X_i X_j]$ ,  $\mathbb{E}[X_i X_j X_k]$ , etc. Hence, marginal-only criteria cannot in general control the joint distribution. We next show that the CF and its derivatives at the origin provide an exact, layered access to all mixed moments/cumulants, thereby capturing dependence.

### 8.3 Taylor Coefficients and Mixed Moments/ Cumulants

**Proposition 10** (Derivative–Moment Correspondence). *Suppose  $\mathbb{E}|X^\alpha| < \infty$ . Then*

$$\partial^\alpha \varphi(0) = i^{|\alpha|} \mathbb{E}[X^\alpha]$$

for every multi-index  $\alpha \in \mathbb{N}_0^d$ .

*Proof.* Write  $e^{i\omega^\top X} = \prod_{j=1}^d e^{i\omega_j X_j}$  and differentiate. Each differentiation w.r.t.  $\omega_j$  contributes a factor  $iX_j$ . Thus  $\partial^\alpha e^{i\omega^\top X} = (iX_1)^{\alpha_1} \dots (iX_d)^{\alpha_d} e^{i\omega^\top X}$ . Taking expectation and setting  $\omega = 0$  yields the claim, with the interchange of differentiation and expectation justified by  $\mathbb{E}|X^\alpha| < \infty$ .  $\square$

**Low-order examples.** For  $\alpha = e_j$  (the  $j$ -th unit vector),  $\partial^\alpha \varphi(0) = i\mathbb{E}[X_j]$  (means). For  $\alpha = e_j + e_k$ ,  $\partial^\alpha \varphi(0) = -\mathbb{E}[X_j X_k]$  (mixed second moments). For  $\alpha = e_i + e_j + e_k$ ,  $\partial^\alpha \varphi(0) = -i\mathbb{E}[X_i X_j X_k]$  (triple mixed moments).

**Proposition 11** (Multivariate Taylor expansion of  $\varphi$  at 0). *Under  $\mathbb{E}|X^\alpha| < \infty$  for all  $|\alpha| \leq m$ ,*

$$\varphi(\omega) = \sum_{|\alpha| \leq m} \frac{1}{\alpha!} \partial^\alpha \varphi(0) \omega^\alpha + o(\|\omega\|^m) = \sum_{|\alpha| \leq m} \frac{i^{|\alpha|}}{\alpha!} \mathbb{E}[X^\alpha] \omega^\alpha + o(\|\omega\|^m). \quad (10)$$

Define  $\psi(\omega) = \log \varphi(\omega)$  (well-defined near 0 since  $\varphi(0) = 1$  and  $\varphi$  is continuous).

**Definition 1** (Log-Characteristic Function and Mixed Cumulants). *For a multi-index  $\alpha \neq 0$  with  $\mathbb{E}|X^\alpha| < \infty$ , the mixed cumulant is*

$$\kappa_\alpha := i^{-|\alpha|} \partial^\alpha \psi(0)$$

(with  $\kappa_0 := 0$  by convention).

**Proposition 12** (Derivative–Cumulant correspondence). *Assume the moments needed are finite so that  $\psi$  is  $|\alpha|$ -times differentiable at 0. Then the coefficients of the multivariate Taylor series of  $\psi$  at 0 equal the mixed cumulants:*

$$\psi(\omega) = \sum_{|\alpha| \geq 1} \frac{1}{\alpha!} \partial^\alpha \psi(0) \omega^\alpha = \sum_{|\alpha| \geq 1} \frac{i^{|\alpha|}}{\alpha!} \kappa_\alpha \omega^\alpha \quad (\text{convergent near } 0). \quad (11)$$

*Proof.* By the chain rule,  $\nabla \psi = \nabla \varphi / \varphi$  and  $\nabla^2 \psi = (\nabla^2 \varphi) / \varphi - (\nabla \varphi \nabla \varphi^\top) / \varphi^2$ , etc. Evaluating at 0 and using Prop. 10:

$$\begin{aligned} \nabla \psi(0) &= i \mathbb{E}[X], \\ \nabla^2 \psi(0) &= -(\mathbb{E}[XX^\top] - \mathbb{E}[X]\mathbb{E}[X]^\top) = -\text{Cov}(X). \end{aligned} \quad (12)$$

Higher-order derivatives of  $\psi$  yield the classical cumulant tensors (via the multivariate Faà di Bruno formula). Collecting terms gives the stated Taylor expansion with coefficients  $\partial^\alpha \psi(0) = i^{|\alpha|} \kappa_\alpha$ .  $\square$

**Interpretation and examples.**  $\kappa_{e_j} = \mathbb{E}[X_j]$  (means);  $\kappa_{e_j+e_k} = \text{Cov}(X_j, X_k)$  (covariance); third-order  $\kappa_{e_i+e_j+e_k}$  measure non-Gaussian triple interactions (synergy/redundancy). In general, *cross-block independence* forces all mixed cumulants spanning the blocks to vanish, making cumulants a clean diagnostic of dependence.

## Implications

- Moments:  $\partial^\alpha \varphi(0) = i^{|\alpha|} \mathbb{E}[X^\alpha]$  reveals *all* mixed moments at each order  $|\alpha|$ .
- Cumulants:  $\partial^\alpha \psi(0) = i^{|\alpha|} \kappa_\alpha$  isolates genuine interactions (they are additive for independent sums and vanish across independent groups).
- Hence, matching  $\varphi$  (or  $\psi$ ) near 0 across *all* multi-indices matches *all* mixed moments/cumulants, aligning dependence at every order.

## 9 Additional Analysis on PD-CFD

Using the CUB-200-2011.<sup>1</sup> bird dataset illustrated in Fig. 3, we evaluate how PD-CFD improves the phase-focusing behavior of NCFM.<sup>2</sup> The CUB-200-2011 dataset contains extremely fine-grained categories where discriminative information resides mainly in high-frequency components such as feather textures. Therefore, the loss must be sensitive to detailed structural variations.

We first revisit the original NCFM loss:

$$\text{Chf}(t; f) = \alpha \left( \left| \Phi_{f(x)}(t) - \Phi_{\hat{f}(x)}(t) \right|^2 \right) + (1 - \alpha) \left| \Phi_{f(x)}(t) \right| \left| \Phi_{\hat{f}(x)}(t) \right| \left( 1 - \cos \left( a_{f(x)}(t) - a_{\hat{f}(x)}(t) \right) \right) \quad (13)$$

where  $\Phi$  denotes the characteristic function magnitude and  $a(\cdot)$  the phase.

Although NCFM applies different hyperparameters to amplitude and phase, the phase term is still tightly entangled with amplitude through the multiplicative magnitude factor. As frequency increases, the amplitude decays, causing the phase contribution to be increasingly suppressed. As a result, the mid- and high-frequency phases—even those

<sup>1</sup>C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, *Caltech-UCSD Birds-200-2011 Dataset*, Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

<sup>2</sup>Wang et al., *Dataset Distillation with Neural Characteristic Function: A Minmax Perspective*, in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2025, pp. 25570–25580.

that still contain meaningful structural information—are overwhelmed and incorrectly treated as noise. To illustrate this behavior, we visualize phase differences across frequencies in Fig.1: low-frequency phase remains stable, the mid-frequency phase begins to drift, and the high-frequency phase becomes dominated by boundary noise. However, some med- and med-high-frequency phase components are still reliable and should not be discarded prematurely.

To restore these informative phase components, we extend the NCFM loss by adding an explicit phase constraint defined for each sampled frequency  $\omega \sim p(t)$  that is decoupled from amplitude :

$$\begin{aligned} \text{Chf}(t; f) = & \alpha \left( \left| \Phi_{f(x)}(t) - \Phi_{\hat{f}(x)}(t) \right|^2 \right) \\ & + (1 - \alpha) \left| \Phi_{f(x)}(t) \right| \left| \Phi_{\hat{f}(x)}(t) \right| \left( 1 - \cos \left( a_{f(x)}(t) - a_{\hat{f}(x)}(t) \right) \right) \\ & + \frac{\lambda_p}{1 + \beta \|\omega\|^2} \left( \theta_{f(x)}(t) - \theta_{\hat{f}(x)}(t) \right)^2, \end{aligned} \quad (14)$$

This additional phase term explicitly extracts the remaining reliable phase information that is otherwise buried by amplitude attenuation in the original NCFM formulation. In particular, it allows mid-frequency phase—which is still semantically meaningful but degraded by the amplitude coupling—to be effectively preserved.

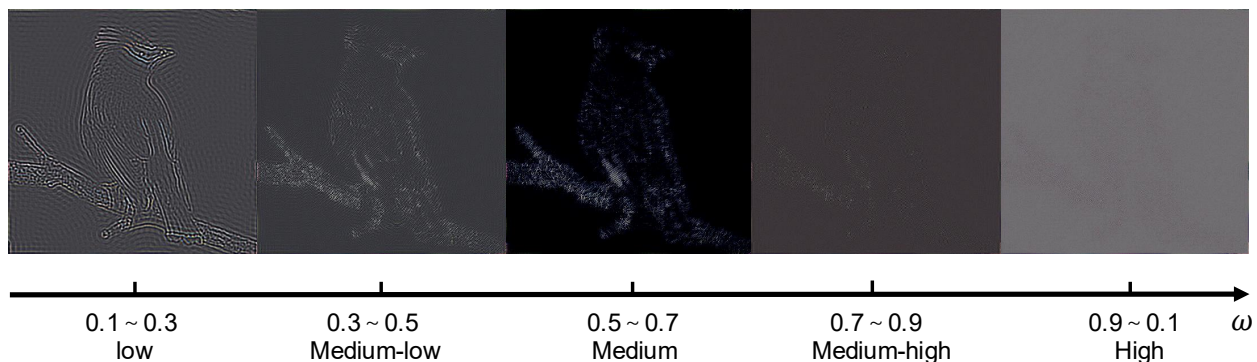


Figure 1: Phase information under different frequency band



Figure 2: Comparison between generated data with original loss function and PD-CFD loss.

We verify the effect of this decoupled phase constraint on downstream classification accuracy, and the results are illustrated in Table. 8 (see Fig. 2 for an example of the synthetic image). The results match our expectation: when  $\lambda_p = 0.3$ , the added phase regularization significantly strengthens high-frequency alignment, especially for the feather and edge regions in the CUB images.

A similar phenomenon is observed on the DTD and RESISC-45 datasets(illustrated in Fig. 4 Fig. 5 respectively). The DTD texture dataset contains a large amount of rapidly oscillating material patterns characterized by abrupt edge transitions, fine-scale boundaries, and dense local contour variations. Likewise, the RESISC-45 remote-sensing dataset

Table 8: Effect of phase regularization  $\lambda_p$  on accuracy.

$\lambda_p$	Accuracy (%)	Improvement
0.0	23.80	0
0.1	24.15	1.47%
0.2	25.20	5.88%
0.3	28.35	19.12%

exhibits substantial geometric complexity: land-cover boundaries, building edges, road networks, rooftop outlines, and other high-frequency landscape structures that vary sharply across classes. These discriminative details are dominated by high-order moment information in the frequency domain. Traditional metrics such as MSE or CE are inherently insensitive to such structural discontinuities—they primarily respond to low-order statistics and smooth variations, and therefore fail to capture the fine-grained distributional differences arising from high-frequency contours and edge transitions.

In contrast, our PD-CFD formulation introduced earlier can naturally recover these informative distributional discrepancies. By removing the amplitude-induced suppression of mid- and high-frequency phase, PD-CFD preserves the remaining reliable phase components that encode exactly these texture- and edge-based variations. Consequently, FAST maintains strong performance even on challenging datasets such as DTD and RESISC-45, where structural information is dominated by, high-frequency features and where competing baselines degrade substantially.

## 10 Experiments on Language Tasks

### 10.1 Setup

To further evaluate the generalization capability of FAST on LLM datasets, we conduct experiments on the Alpaca instruction-following dataset. Following our main setup, we adopt LLaMA2-7B as the backbone and apply LoRA-based fine-tuning using core-sets sampled at retention rates of 10%, 20%, and 30%. For evaluation, we report performance on the MMLU benchmark, which spans 57 diverse subjects and serves as a rigorous protocol for assessing broad knowledge and reasoning abilities. We compare FAST against the NMS baseline<sup>3</sup>, the current SOTA method, under identical experimental settings.

### 10.2 Results

The MMLU results are shown in Fig. 6. Across all keep ratios, FAST consistently outperforms the NMS baseline. While NMS attains an average accuracy of approximately 33%, FAST increases the performance to 39%, corresponding to a relative improvement of about 18%. These results indicate that FAST is highly effective at selecting and preserving semantically informative instances. Moreover, the strong performance gains suggest that high-level semantic structure can be retained through spectral-graph-based distribution alignment alone, without explicit reliance on neural feature extractors. By leveraging geometric and frequency-domain signals, FAST preserves the underlying semantic neighborhood relations required for downstream reasoning. This demonstrates the robustness and broad generalization ability of FAST, even on tasks that rely heavily on semantic consistency.

<sup>3</sup>Boran Zhao et al., *NMS: Efficient Edge DNN Training via Near-Memory Sampling on Manifolds*, arXiv preprint arXiv:2508.02313, 2025.



Figure 3: CUB-200-2011 Bird Dataset

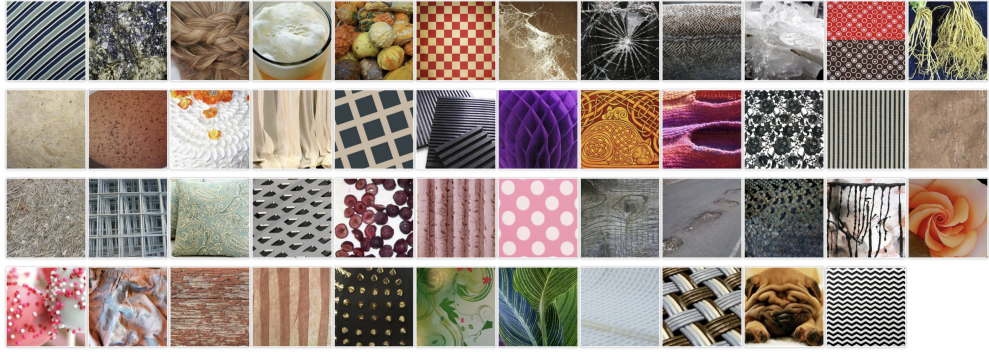


Figure 4: DTD Texture Dataset

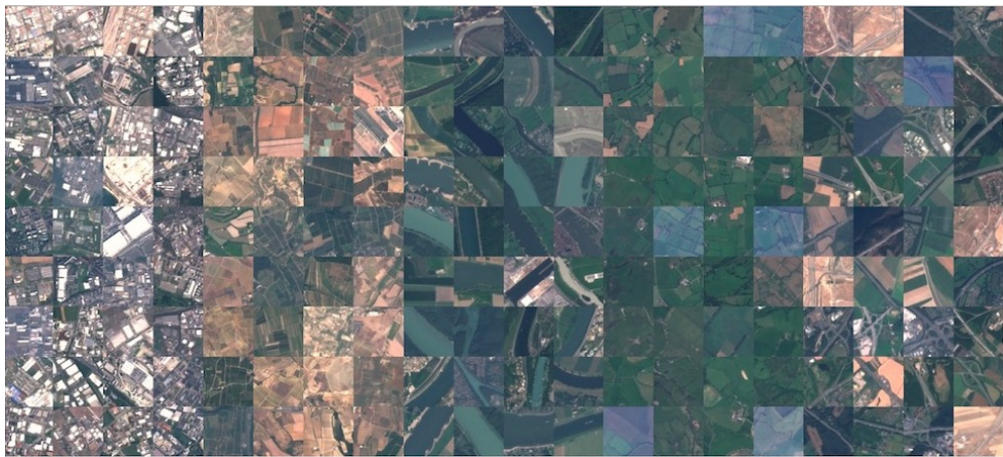
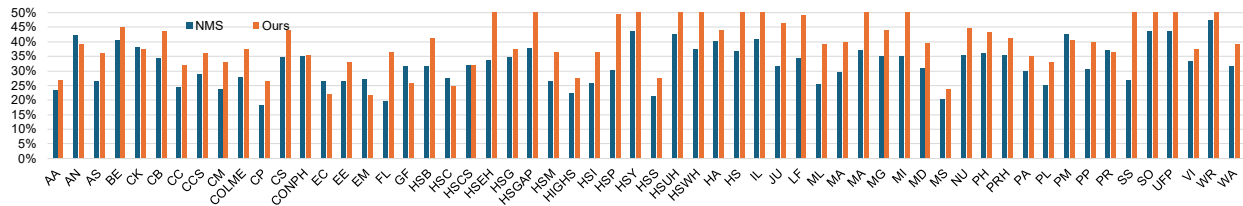
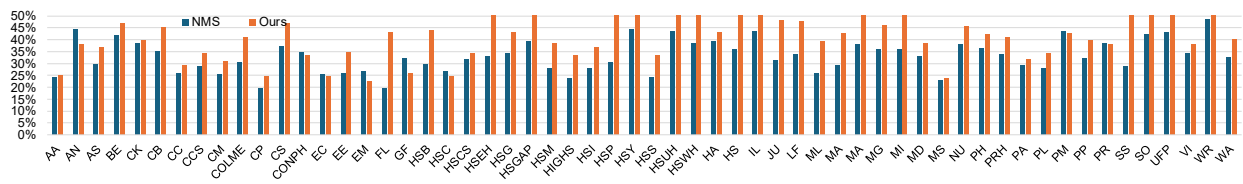


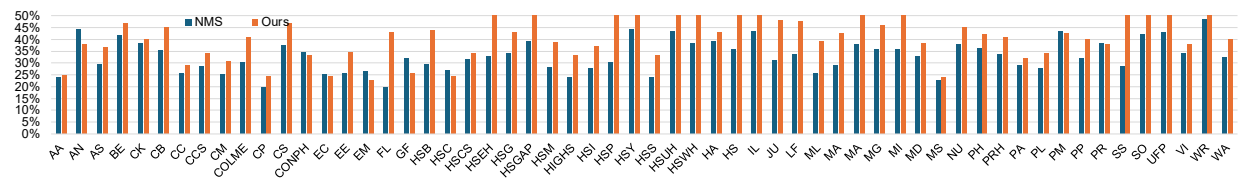
Figure 5: RESISC-45 Remote Sensing Dataset



(a) Keep ratio 10%



(b) Keep ratio 20%



(c) Keep ratio 30%

Figure 6: MMLU accuracy comparison at different keep ratios.