

APPENDIX: MATHEMATICAL THEOREMS AND COROLLARY PROOFS

Theorem 1 For any two data distributions \mathcal{P}_1 and \mathcal{P}_2 , if their distribution difference in the original space is $\mathcal{D}(\mathcal{P}_1, \mathcal{P}_2) < \varepsilon$, then there exists an optimal set of diffusion parameters $\{\alpha_t^*\}_{t=1}^T$ such that the distribution difference after DMFE feature extraction is $\mathcal{D}(\phi_{DMFE}^*(\mathcal{P}_1), \phi_{DMFE}^*(\mathcal{P}_2)) > \delta$, where $\delta \gg \varepsilon$, and satisfies:

$$\{\alpha_t^*\}_{t=1}^T = \arg \max_{\{\alpha_t\}_{t=1}^T} \mathcal{D}(\phi_{DMFE}(\mathcal{P}_1), \phi_{DMFE}(\mathcal{P}_2)) - \lambda \sum_{t=1}^T (1 - \alpha_t) \quad (1)$$

where λ is a regularization parameter balancing distribution separability and diffusion process stability.

Proof 1 The Wasserstein-1 distance is defined via optimal transport: $W_1(\mathcal{P}, \mathcal{Q}) = \inf_{\gamma \in \Gamma(\mathcal{P}, \mathcal{Q})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} [\|\mathbf{x} - \mathbf{y}\|]$ where $\Gamma(\mathcal{P}, \mathcal{Q})$ is the set of all couplings. For a Lipschitz-continuous function ϕ with constant $L_\phi = \sup_{\mathbf{x} \neq \mathbf{y}} \|\phi(\mathbf{x}) - \phi(\mathbf{y})\| / \|\mathbf{x} - \mathbf{y}\|$, the contraction property gives:

$$W_1(\phi(\mathcal{P}), \phi(\mathcal{Q})) \leq L_\phi \cdot W_1(\mathcal{P}, \mathcal{Q}) \quad (2)$$

For a composition of Lipschitz functions $\phi = \phi_K \circ \dots \circ \phi_1$ with individual constants L_i , the composed Lipschitz constant is the product:

$$L_\phi = \prod_{i=1}^K L_i \quad (3)$$

The DMFE architecture decomposes as sequential operations. For the diffusion process, at each timestep:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \varepsilon_{t-1} \quad (4)$$

The Lipschitz constant of this operation satisfies:

$$L_t = \sup_{\mathbf{x}_{t-1} \neq \mathbf{x}'_{t-1}} \frac{\|\sqrt{\alpha_t} \mathbf{x}_{t-1} - \sqrt{\alpha_t} \mathbf{x}'_{t-1}\|}{\|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|} = \sqrt{\alpha_t} \quad (5)$$

Composing T diffusion steps yields:

$$L_{diff} = \prod_{t=1}^T \sqrt{\alpha_t} = \left(\prod_{t=1}^T \alpha_t \right)^{1/2} \quad (6)$$

For conservative diffusion schedules with $\alpha_t \in [0.9, 0.99]$, we have $\prod_t \alpha_t \in [0.64, 1.00]$, giving $L_{diff} \in [0.80, 1.00]$.

The frequency domain transforms (Fourier \mathcal{F} and wavelet \mathcal{W}) are orthogonal linear operators satisfying:

$$\|\mathcal{F}(\mathbf{x})\| = \|\mathbf{x}\|, \quad \|\mathcal{W}(\mathbf{x})\| = \|\mathbf{x}\| \quad (7)$$

thus $L_{freq} = 1$.

The convolutional feature extractors apply K_{conv} layers. Under spectral normalization, the normalized weight matrix has $\|\mathbf{W}_i^{norm}\|_{op} = 1$. With learnable scaling factor s_i , the actual Lipschitz constant becomes:

$$L_i = (1 + s_i) \|\mathbf{W}_i^{norm}\|_{op} = 1 + s_i \quad (8)$$

Composing all convolutional layers:

$$L_{conv} = \prod_{i=1}^{K_{conv}} (1 + s_i) \quad (9)$$

Under the constraint $\sum_{i=1}^{K_{conv}} s_i \leq S_{max}$, maximizing the product occurs with uniform scaling $s_i = S_{max}/K_{conv}$:

$$L_{conv}^{max} = \left(1 + \frac{S_{max}}{K_{conv}}\right)^{K_{conv}} \quad (10)$$

- 11 For $S_{max} = 1$ and $K_{conv} = 16$: $(1.0625)^{16} \approx 2.64$. For $S_{max} = 1.5$ and $K_{conv} = 12$: $(1.125)^{12} \approx 3.22$.
The covariance pooling and patch aggregation operations contribute:

$$L_{pool} = 1 + c_{pool} \in [1.05, 1.15], \quad L_{fusion} = 1 + c_{fusion} \in [1.05, 1.20] \quad (11)$$

- 12 where c_{pool}, c_{fusion} are empirically determined constants from feature statistics.
By the composition rule, the total Lipschitz constant is:

$$L_{total} = L_{diff} \cdot L_{freq} \cdot L_{conv} \cdot L_{pool} \cdot L_{fusion} \quad (12)$$

For conservative settings with $L_{diff} \approx 0.95$, $L_{conv} \approx 2.64$, $L_{pool} \approx 1.10$, $L_{fusion} \approx 1.10$:

$$L_{total} \approx 0.95 \times 1 \times 2.64 \times 1.10 \times 1.10 \approx 3.0 \quad (13)$$

For aggressive settings with $L_{diff} \approx 1.00$, $L_{conv} \approx 3.22$, $L_{pool} \approx 1.15$, $L_{fusion} \approx 1.20$:

$$L_{total} \approx 1.00 \times 1 \times 3.22 \times 1.15 \times 1.20 \approx 4.45 \quad (14)$$

By the Lipschitz contraction property:

$$W_1(\phi_{DMFE}(\mathcal{P}_1), \phi_{DMFE}(\mathcal{P}_2)) \leq L_{total} \cdot W_1(\mathcal{P}_1, \mathcal{P}_2) = L_{total} \cdot \epsilon_0 \quad (15)$$

- 13 which establishes the stated amplification bound.

- 14 **Corollary 1** Given a set of AI generators $\{G_1, G_2, \dots, G_M\}$ and their corresponding image distributions
15 $\{\mathcal{P}_g^1, \mathcal{P}_g^2, \dots, \mathcal{P}_g^M\}$, if DMFE can effectively distinguish between the real image distribution \mathcal{P}_r and
16 the distribution of any known generator, then for the distribution \mathcal{P}_g^{new} produced by a newly emerging
17 unknown generator G_{new} , DMFE can still maintain a certain discrimination ability, with its generalization
18 performance lower bound:

$$\mathcal{D}(\phi_{DMFE}^*(\mathcal{P}_r), \phi_{DMFE}^*(\mathcal{P}_g^{new})) \geq \min_{j \in \{1, 2, \dots, M\}} \mathcal{D}(\phi_{DMFE}^*(\mathcal{P}_r), \phi_{DMFE}^*(\mathcal{P}_g^j)) - \gamma \cdot \max_{j \in \{1, 2, \dots, M\}} \mathcal{D}(\mathcal{P}_g^j, \mathcal{P}_g^{new}) \quad (16)$$

- 19 where γ is a constant related to the complexity of the DMFE model. This corollary indicates that
20 DMFE's generalization ability to unknown generators depends on its discrimination performance for
21 known generators and the similarity between the distributions of new generators and known generators.

Proof 2 The Wasserstein distance satisfies the triangle inequality and its reverse form. By the reverse triangle inequality:

$$|W_1(\mathcal{P}_1, \mathcal{P}_3) - W_1(\mathcal{P}_1, \mathcal{P}_2)| \leq W_1(\mathcal{P}_2, \mathcal{P}_3) \quad (17)$$

Rearranging:

$$W_1(\mathcal{P}_1, \mathcal{P}_3) \geq W_1(\mathcal{P}_1, \mathcal{P}_2) - W_1(\mathcal{P}_2, \mathcal{P}_3) \quad (18)$$

Apply this in the feature space. For any known generator index j :

$$W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^{new})) \geq W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^j)) - W_1(\phi_{DMFE}(\mathcal{P}_g^j), \phi_{DMFE}(\mathcal{P}_g^{new})) \quad (19)$$

By the Lipschitz contraction property from Theorem ??:

$$W_1(\phi_{DMFE}(\mathcal{P}_g^j), \phi_{DMFE}(\mathcal{P}_g^{new})) \leq L_{total} \cdot W_1(\mathcal{P}_g^j, \mathcal{P}_g^{new}) \leq L_{total} \cdot \epsilon_{novel} \quad (20)$$

Substituting back:

$$W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^{new})) \geq W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^j)) - L_{total} \cdot \epsilon_{novel} \quad (21)$$

Taking the minimum over all j :

$$W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^{new})) \geq \min_j W_1(\phi_{DMFE}(\mathcal{P}_r), \phi_{DMFE}(\mathcal{P}_g^j)) - L_{total} \cdot \epsilon_{novel} \quad (22)$$

22 This demonstrates that discrimination ability on unknown generators degrades linearly with their
23 distributional distance from known generators, with proportionality constant L_{total} .

24 **Theorem 2** Given two data distributions $\phi_{DMFE}^*(\mathcal{P}_r)$ and $\phi_{DMFE}^*(\mathcal{P}_g)$ in the feature space extracted by
25 DMFE, if CAViT can effectively capture higher-order relationships of these features, then the cascaded
26 ensemble model DCAT can further increase the distribution difference, satisfying:

$$\begin{aligned} \mathcal{D}(\phi_{DCAT}(\mathcal{P}_r), \phi_{DCAT}(\mathcal{P}_g)) &\geq \eta \cdot \mathcal{D}(\phi_{DMFE}^*(\mathcal{P}_r), \phi_{DMFE}^*(\mathcal{P}_g)) \\ &\quad + \kappa \cdot KL(\phi_{DMFE}^*(\mathcal{P}_r) \parallel \phi_{DMFE}^*(\mathcal{P}_g)) \\ &\quad + \gamma \cdot (1 - \exp(-\beta \cdot \|\mathbb{E}_{\mathbf{x} \sim \mathcal{P}_r}[\phi_{DMFE}^*(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_g}[\phi_{DMFE}^*(\mathbf{x})]\|_2^2)) \\ &\quad - \alpha \cdot \sqrt{\text{Var}_{\mathbf{x} \sim \mathcal{P}_r}[\phi_{DMFE}^*(\mathbf{x})] + \text{Var}_{\mathbf{x} \sim \mathcal{P}_g}[\phi_{DMFE}^*(\mathbf{x})]} \end{aligned} \quad (23)$$

27 where $\mathcal{D}(\cdot, \cdot)$ is the distribution difference measure, as defined in equation (1); $\eta > 1$ is the basic
28 enhancement factor, positively correlated with CAViT's expressive power and parameter count; $\kappa > 0$ is
29 the weight coefficient for KL divergence; $KL(\cdot \parallel \cdot)$ is the Kullback-Leibler divergence; $\gamma > 0$ and $\beta > 0$ are
30 coefficients for the mean difference term; $\alpha > 0$ is the variance penalty coefficient; \mathbb{E} and Var represent
31 expectation and variance, respectively; $\phi_{DCAT}(\mathcal{P})$ represents the distribution after mapping distribution
32 \mathcal{P} through the DCAT model. This theorem shows that CAViT not only linearly amplifies the distribution
33 differences produced by DMFE but also further enhances discrimination ability through modeling of
34 higher-order statistics.

Proof 3 The global path in CAViT applies L_g sequential attention layers. The softmax operation in
attention has Lipschitz constant 1 because the Jacobian eigenvalues of softmax satisfy $\lambda_{softmax} \leq 1$.
The query-key bilinear form $\mathbf{QW}^Q(\mathbf{KW}^K)^T$ has operator norm at most 1 under spectral normalization.
Each attention head produces dimension $d_v = D/h$. When h independent attention pathways merge
through output projection \mathbf{W}^O , gradient amplification proportional to \sqrt{h} occurs (from variance of sum).
Simultaneously, reduced per-head dimension $d_k = D/h$ increases attention logits by factor \sqrt{h} relative to
full dimension, amplifying gradient flow. The composite effect yields:

$$L_{attn} = 1 + c_1 \frac{\sqrt{h}}{\sqrt{d_k}} = 1 + c_1 \frac{h}{\sqrt{D}} \quad (24)$$

Composing L_g such layers via the product rule:

$$L_{global} = \prod_{l=1}^{L_g} (1 + \epsilon_{attn}) = (1 + \epsilon_{attn})^{L_g} =: \eta \quad (25)$$

For typical settings with $L_g = 12$, $h = 8$, $D = 512$: $\epsilon_{attn} = 0.5 \times 8/\sqrt{512} \approx 0.088$, giving $\eta =$
(1.088)¹² ≈ 2.82 . Thus:

$$W_1(\text{GlobalPath}(\mathcal{P}_r), \text{GlobalPath}(\mathcal{P}_g)) \geq \eta \cdot \delta_0 \quad (26)$$

The KL divergence between distributions is:

$$KL(\mathcal{P}_r \parallel \mathcal{P}_g) = \int \log \frac{p_r(\mathbf{x})}{p_g(\mathbf{x})} p_r(\mathbf{x}) d\mathbf{x} \quad (27)$$

Pinsker's inequality relates KL divergence to total variation distance:

$$TV(\mathcal{P}_r, \mathcal{P}_g) = \frac{1}{2} \int |p_r(\mathbf{x}) - p_g(\mathbf{x})| d\mathbf{x} \leq \sqrt{\frac{KL(\mathcal{P}_r \| \mathcal{P}_g)}{2}} \quad (28)$$

For bounded supports with diameter D_{supp} , Wasserstein distance relates to total variation by:

$$W_1(\mathcal{P}_r, \mathcal{P}_g) \geq \frac{1}{D_{supp}} TV(\mathcal{P}_r, \mathcal{P}_g) \geq C_{supp} \sqrt{KL(\mathcal{P}_r \| \mathcal{P}_g)} \quad (29)$$

where $C_{supp} = 1/D_{supp}$. The local path's graph attention network with Lipschitz constant κ extracts this information:

$$W_1(\text{LocalPath}(\mathcal{P}_r), \text{LocalPath}(\mathcal{P}_g)) \geq \kappa \cdot C_{supp} \sqrt{KL(\mathcal{P}_r \| \mathcal{P}_g)} \quad (30)$$

The patch relationship modeling computes pairwise interactions with Lipschitz constant $L_{patch} = \gamma$. By Kantorovich duality:

$$W_1(\mathcal{P}_1, \mathcal{P}_2) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathcal{P}_1}[f] - \mathbb{E}_{\mathcal{P}_2}[f] \quad (31)$$

Choosing linear functional in direction of $\Delta\mu_0$:

$$W_1(\mathcal{P}_1, \mathcal{P}_2) \geq \Delta\mu_0 - \sqrt{2\sigma_0^2} \quad (32)$$

For $\Delta\mu_0 \gg \sigma_0$, this approximates:

$$W_1(\mathcal{P}_1, \mathcal{P}_2) \geq \Delta\mu_0 \left(1 - \exp\left(-\frac{(\Delta\mu_0)^2}{c_{mv}\sigma_0^2}\right) \right) \quad (33)$$

where $c_{mv} \approx 4$ for Gaussian distributions. The patch modeling contributes:

$$W_1(\text{PatchPath}(\mathcal{P}_r), \text{PatchPath}(\mathcal{P}_g)) \geq \gamma \left(1 - \exp\left(-\frac{(\Delta\mu_0)^2}{c_{mv}\sigma_0^2}\right) \right) \quad (34)$$

The variance penalty $-\alpha\sigma_0$ prevents overfitting to high-overlap regions. The four mechanisms—spectral amplification, KL divergence, mean separation, and stability—operate on distinct statistical aspects. Their contributions can be combined additively via Kantorovich duality, as different optimal Lipschitz functions capture different distributional aspects:

$$W_1(\phi_{DCAT}(\mathcal{P}_r), \phi_{DCAT}(\mathcal{P}_g)) \geq \eta\delta_0 + \kappa C_{supp} \sqrt{KL} + \gamma \left(1 - \exp\left(-\frac{(\Delta\mu_0)^2}{c_{mv}\sigma_0^2}\right) \right) - \alpha\sigma_0 \quad (35)$$

35 **Corollary 2** Based on Theorem ??, for any distribution pair satisfying $\mathcal{D}(\mathcal{P}_r, \mathcal{P}_g) < \varepsilon$, the discrimi-
 36 nation performance of the DCAT ensemble model under optimal parameter configuration has a lower
 37 bound:

$$\begin{aligned} AUC(\phi_{DCAT}) &\geq 1 - \exp\left(-\frac{\eta \cdot \delta^2}{2}\right) \\ &\quad - \exp(-\lambda_1 \cdot KL(\mathcal{P}_r \| \mathcal{P}_g) - \lambda_2 \cdot JS(\mathcal{P}_r, \mathcal{P}_g)) \\ &\quad - \frac{C_1}{\sqrt{n}} \cdot \left(1 + \sqrt{\log \frac{1}{\rho}} \right) \\ &\quad + C_2 \cdot \left(1 - \frac{1}{1 + \exp(\omega \cdot (\|\bullet\|_F - \tau))} \right) \end{aligned} \quad (36)$$

38 where $AUC(\phi_{DCAT})$ represents the area under the receiver operating characteristic curve of model
 39 ϕ_{DCAT} on the binary classification task of real versus AI-generated images; $\delta > 0$ is the minimum dis-
 40 tribution difference achieved by DMFE in the feature space, as guaranteed by Theorem ??; $\eta > 1$ is

41 the enhancement factor provided by CAViT; λ_1 and λ_2 are weight coefficients for information-theoretic
 42 measures; $KL(\cdot\|\cdot)$ is the Kullback-Leibler divergence; $JS(\cdot, \cdot)$ is the Jensen-Shannon divergence; C_1
 43 is a constant related to sample complexity; n is the number of training samples; ρ is a confidence
 44 parameter; C_2 , ω , and τ are parameters related to model complexity; $\|\bullet\|_F$ is the Frobenius norm of the
 45 model parameters. This corollary comprehensively considers the influence of distribution differences,
 46 information-theoretic measures, sample complexity, and model complexity on model performance, ensuring
 47 that DCAT can achieve effective discrimination even when real and AI-generated images are extremely
 48 similar in the original image space, with performance improving as η and δ increase.

Proof 4 For sub-Gaussian distributions in \mathbb{R} with means μ_r, μ_g and combined variance σ_{DCAT}^2 , the optimal AUC for the likelihood ratio test satisfies:

$$AUC_{opt} = P(Z_g > Z_r) \approx \Phi\left(\frac{\mu_g - \mu_r}{\sqrt{2}\sigma_{DCAT}}\right) \quad (37)$$

where Φ is the standard normal CDF. Using the asymptotic expansion $1 - \Phi(x) \approx (1/\sqrt{2\pi}x) \exp(-x^2/2)$:

$$AUC_{opt} \geq 1 - \frac{1}{2} \exp\left(-\frac{(\mu_g - \mu_r)^2}{8\sigma_{DCAT}^2}\right) \quad (38)$$

The Wasserstein distance between output distributions equals the mean difference:

$$W_1(Z_r, Z_g) = |\mu_g - \mu_r| = \Delta\mu_{DCAT} \quad (39)$$

By Theorem ??, this separation is at least $\eta\delta_0$:

$$\Delta\mu_{DCAT} \geq \eta\delta_0 \quad (40)$$

Therefore:

$$AUC_{empirical} \geq 1 - \frac{1}{2} \exp\left(-\frac{(\eta\delta_0)^2}{8\sigma_{DCAT}^2}\right) \quad (41)$$

The empirical AUC on finite training data differs from true AUC. By Vapnik-Chervonenkis theory, for a hypothesis class with VC dimension V_{VC} :

$$|AUC_{empirical} - AUC_{true}| \leq \frac{C_1}{\sqrt{n}} \left(1 + \sqrt{\log(1/\rho)}\right) \quad (42)$$

where $C_1 = O(\sqrt{V_{VC}})$ and $V_{VC} = O(LD \log D)$ for deep networks with L layers of width D . Thus:

$$AUC_{true} \geq AUC_{empirical} - \frac{C_1}{\sqrt{n}} \left(1 + \sqrt{\log(1/\rho)}\right) \quad (43)$$

The KL divergence contributes through Chernoff information theory. The error probability under the optimal classifier satisfies:

$$P_{error} \leq \exp(-C_{Chernoff} \cdot KL(\mathcal{P}_r\|\mathcal{P}_g)) \quad (44)$$

for Chernoff constant $C_{Chernoff} > 0$. This translates to AUC contribution:

$$AUC_{info} \geq 1 - \exp(-\lambda_1 KL(\mathcal{P}_r\|\mathcal{P}_g)) \quad (45)$$

Combining all three sources—Wasserstein separation (architecture), generalization (learning theory), and information-theoretic advantage (distribution properties):

$$AUC(\phi_{DCAT}) \geq \left(1 - \frac{1}{2} \exp\left(-\frac{(\eta\delta_0)^2}{8\sigma_{DCAT}^2}\right)\right) - \frac{C_1}{\sqrt{n}} \left(1 + \sqrt{\log(1/\rho)}\right) + (1 - \exp(-\lambda_1 KL(\mathcal{P}_r\|\mathcal{P}_g))) \quad (46)$$