#### A. Appendix

#### **B.** Statistical Interpretation of the Entropy-based Scores

In this section, we aim to statistically interpret the defined conditional diversity scores as the expectation of prompt-specific entropy H(X|T=t). First, we derive the statistic estimated from empirical samples by the entropy-based scores, and then we connect the conditional entropy measure to the expectation of unconditional entropy values. According to the Schur product theorem, the Hadamard product  $K_X \odot K_T$  of PSD kernel matrices  $K_X$ ,  $K_T$  will also be a PSD kernel matrix. We note that the kernel matrix corresponds to the following feature map  $\phi_{X,T}: \mathcal{X} \times \mathcal{T} \to \mathbb{R}^{d_x d_t}$  where  $\otimes$  denotes the Kronecker product:

$$\phi_{X,T}([x,t]) = \phi_X(x) \otimes \phi_T(t)$$

The above holds due to the identity  $\langle \phi_{X,T}([x,t]), \phi_{X,T}([x',t']) \rangle = k_X(x,x')k_T(t,t')$ . The following proposition formulates kernel-based conditional entropy and mutual information using  $\phi_{X,T}$ .

**Proposition 1.** Consider the kernel matrices  $K_X$  for samples  $x_1, \ldots, x_n$  and  $K_T$  for samples  $t_1, \ldots, t_n$ . Then,  $\frac{1}{n}K_X \odot K_T$  (used for defining joint entropy  $H_{\alpha}(X,T)$ ) share the same non-zero eigenvalues with the following kernel covariance matrix:

$$\widehat{C}_{X,T} := \frac{1}{n} \sum_{i=1}^{n} \phi_{X,T} ([x_i, t_i]) \phi_{X,T} ([x_i, t_i])^{\top} = \frac{1}{n} \sum_{i=1}^{n} \left[ \phi_X (x_i) \phi_X (x_i)^{\top} \right] \otimes \left[ \phi_T (t_i) \phi_T (t_i)^{\top} \right]$$

*Proof.* We defer the proof to the Appendix C.1.

**Corollary 1.** Consider the composite feature map  $\phi_{X,T}$  and joint kernel covariance matrix  $\widehat{C}_{X,T}$  defined above. Then, given marginal kernel covariance matrices  $\widehat{C}_X = \frac{1}{n} \sum_{i=1}^n \phi_X(x_i) \phi_X(x_i)^{\top}$ ,  $\widehat{C}_T = \frac{1}{n} \sum_{i=1}^n \phi_T(t_i) \phi_T(t_i)^{\top}$ , the following holds for the defined conditional entropy and mutual information:

$$H_{\alpha}(X|T) = H_{\alpha}(\widehat{C}_{X,T}) - H_{\alpha}(\widehat{C}_{T}), \quad I_{\alpha}(X;T) = H_{\alpha}(\widehat{C}_{X}) + H_{\alpha}(\widehat{C}_{T}) - H_{\alpha}(\widehat{C}_{X,T})$$

*Proof.* The proof is deferred to the Appendix C.2.

Corollary 1 shows that given the underlying covariance matrices  $C_X = \mathbb{E}_{x \sim P_X} \left[ \phi_X(x) \phi_X(x)^\top \right], C_T = \mathbb{E}_{t \sim P_T} \left[ \phi_T(t) \phi_T(t)^\top \right],$  and  $C_{X,T} = \mathbb{E}_{(x,t) \sim P_{X,T}} \left[ \phi_{X,T}([x,t]) \phi_{X,T}([x,t])^\top \right],$  the defined entropy-based scores converge to the following statistics when the sample size n tends to infinity:

$$\widetilde{H}_{\alpha}(X|T) = H_{\alpha}(C_{X,T}) - H_{\alpha}(C_T), \quad \widetilde{I}_{\alpha}(X;T) = H_{\alpha}(C_X) + H_{\alpha}(C_T) - H_{\alpha}(C_{X,T}).$$

Note that the entropy-based statistic  $H_{\alpha}(C_X)$  represents the statistic estimated by the logarithm of the Vendi score defined in [11].

Next, we prove that for a mixture text distribution  $P_T$  where the text variable follows random mode  $G \in \{1, ..., m\}$ , the defined conditional entropy score aggregates the expectation of the unconditional entropy score H(X|G=i) over the m text modes 1, ..., m.

**Theorem 1.** Consider the Gaussian kernel with bandwidth  $\sigma$ . Suppose T follows a mixture distribution  $\sum_{i=1}^m \omega_i P_{T,i}$  where  $\omega_i$  denotes the weight of the ith component  $P_{T,i}$  with mean vector  $\mu_i$  and total variance  $\mathbb{E}_{T \sim P_{T,i}}[\|T - \mu_i\|_2^2] = \sigma_i^2$ . Given the aggregation map  $f(z) = \exp((1 - \alpha)z)$ , for every order  $\alpha \geq 2$ , the matrix-based order- $\alpha$  conditional entropy satisfies the following inequality where  $g(z) = \frac{\alpha}{\alpha - 1} \log \left(\frac{1}{1 - z/\|\omega\|_{\alpha}}\right)$  is an increasing scalar function with g(0) = 0:

$$\left| \widetilde{H}_{\alpha}(X|T) - f^{-1} \left( \mathbb{E}_{I \sim \boldsymbol{\omega}^{\alpha}} \left[ f \left( \widetilde{H}_{\alpha}(X|G=I) \right) \right] \right) \right| \leq 2 g \left( 32 \sum_{i=1}^{k} \omega_{i} \left[ \frac{\sigma_{i}^{2}}{\sigma^{2}} + \sum_{j=1}^{i-1} \exp \left( \frac{-\|\mu_{i} - \mu_{j}\|_{2}^{2}}{\sigma^{2}} \right) \right] \right)$$

*Proof.* The proof is deferred to the Appendix C.3.

The above theorem shows that if the text samples come from m distinct modes satisfying  $\frac{\|\mu_i - \mu_j\|_2}{\sigma} \gg 1$  for every  $i \neq j$  and  $\frac{\mathbb{E}_{P_{T,i}}[\|T - \mu_i\|_2^2]}{\sigma^2} \ll 1$ , then the defined conditional entropy score H(X|T) aggregates the unconditional entropy score H(X|G=i) given the prompt group. Therefore, this result extends the expectation-based interpretation of Shannon conditional entropy to the matrix-based conditional entropy defined in [12].

Based on Theorem 1, we propose a text-type-based diversity evaluation, where we restrict the evaluation of the prompt-based generative model to the prompts in the same group, i.e. the same mode in the mixture text distribution. To do this, we find the eigendirections corresponding to the text clusters by performing an eigendecomposition of text kernel matrix  $\frac{1}{n}K_T = \sum_{i=1}^n \lambda_i v_i v_i^{\mathsf{T}}$  where  $\lambda_1 \geq \cdots \lambda_n$  are the sorted eigenvalues and  $v_1, \ldots, v_n$  are the sorted eigenvectors. Then, we note that the Hadamard product  $\frac{1}{n}K_X \odot K_T$  used for joint entropy can be decomposed as:  $\frac{1}{n}K_X \odot K_T = \sum_{i=1}^n \lambda_i \left(K_X \odot v_i v_i^{\mathsf{T}}\right)$ . As a result, to evaluate the diversity of the generation model, we apply eigendocomposition to each  $K_X \odot v_i v_i^{\mathsf{T}}$ , where  $v_i$  marks the samples in group i, and find the sample indices with the maximum entries on the principal eigenvectors of  $K_X \odot v_i v_i^{\mathsf{T}}$ .

#### C. Proofs

#### C.1. Proof of Proposition 1

First, we observe that for every  $x, x' \in \mathcal{X}$  and  $t, t' \in \mathcal{T}$ , the following holds:

$$\phi_{X,T}([x,t])^{\top}\phi_{X,T}([x',t']) = (\phi_X(x) \otimes \phi_T(t))^{\top} (\phi_X(x') \otimes \phi_T(t'))$$

$$= (\phi_X(x)^{\top}\phi_X(x')) \otimes (\phi_T(t)^{\top}\phi_T(t'))$$

$$= k_X(x,x') \otimes k_T(t,t')$$

$$= k_X(x,x')k_T(t,t')$$

Therefore, the Hadamard product of kernel matrices  $\frac{1}{n}K_X \odot K_T$  can be written as

$$\frac{1}{n}K_X \odot K_T = \frac{1}{n}\Phi_{X,T}\Phi_{X,T}^{\top}$$

in terms of the matrix of samples' feature maps  $\Phi_{X,T} \in \mathbb{R}^{n \times d_x d_t}$  with its ith row being  $\phi_{X,T}([x_i,t_i])$ . We observe that the matrices  $\frac{1}{n}\Phi_{X,T}\Phi_{X,T}^{\top}$  and  $\frac{1}{n}\Phi_{X,T}^{\top}\Phi_{X,T}$  share the same non-zero eigenvalues, that are the square of the singular values of  $\Phi_{X,T}$ . Therefore,  $\frac{1}{n}K_X \odot K_T$  has the same non-zero eigenvalues as the following matrix

$$\frac{1}{n} \Phi_{X,T}^{\top} \Phi_{X,T} = \frac{1}{n} \sum_{i=1}^{n} \phi_{X,T}([x_i, t_i]) \phi_{X,T}([x_i, t_i])^{\top}$$

which is the defined matrix  $\hat{C}_{X,T}$ . Therefore, the proof of the proposition is complete.

#### C.2. Proof of Corollary 1

As we showed in Proposition 1, the Hadarmard product  $\frac{1}{n}K_X \odot K_T$  shares the same non-zero eigenvalues with  $\widehat{C}_{X,T}$ . Also, as noted by [17],  $\frac{1}{n}K_X$  and  $\frac{1}{n}K_T$  have the same non-zero eigenvalues as of  $\widehat{C}_X$  and  $\widehat{C}_T$ , respectively. Since the order- $\alpha$  matrix-based entropy is only a function of the input matrix's non-zero eigenvalues (zero eigenvalues have no impact on the entropy value), we can conclude that

$$H_{\alpha}(X|T) := H_{\alpha}\left(\frac{1}{n}K_X \odot K_T\right) - H_{\alpha}\left(\frac{1}{n}K_T\right)$$
$$= H_{\alpha}(\widehat{C}_{X,T}) - H_{\alpha}(\widehat{C}_T),$$

and also

$$I_{\alpha}(X;T) := H_{\alpha}\left(\frac{1}{n}K_{X}\right) + H_{\alpha}\left(\frac{1}{n}K_{T}\right) - H_{\alpha}\left(\frac{1}{n}K_{X} \odot K_{T}\right)$$
$$= H_{\alpha}(\widehat{C}_{X}) + H_{\alpha}(\widehat{C}_{T}) - H_{\alpha}(\widehat{C}_{X,T}).$$

#### C.3. Proof of Theorem 1

To prove Theorem 1, we begin by showing the following lemma.

**Lemma 1.** Suppose that the kernel function k and variable T satisfy the assumptions in Theorem 1. Then, the following Frobenius norm bound holds for  $C_i = \mathbb{E}[\phi_X(x)\phi_X(x)^\top|G=i]$  where  $G \in \{1, ..., m\}$  is the cluster random variable for text T:

$$\left\| C_{X \otimes T} - \sum_{i=1}^{m} \omega_i C_i \otimes \phi_T(\mu_i) \phi_T(\mu_i)^\top \right\|_F^2 \leq \frac{\sum_{i=1}^{m} 2\omega_i \sigma_i^2}{\sigma^2}.$$

*Proof.* To show this lemma, we define  $T_i$  as a variable distributed as  $P_{T|G=i}$ . Then,

$$\begin{split} & \left\| C_{X \otimes T} - \sum_{i=1}^{m} \omega_{i} C_{i} \otimes \phi(\mu_{i}) \phi(\mu_{i})^{\top} \right\|_{F}^{2} \\ &= \left\| \mathbb{E} \left[ \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \phi_{T}(t) \phi_{T}(t)^{\top} \right] - \sum_{i=1}^{m} \omega_{i} C_{i} \otimes \phi(\mu_{i}) \phi(\mu_{i})^{\top} \right\|_{F}^{2} \\ &= \left\| \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \phi_{T}(t) \phi_{T}(t)^{\top} \middle| G = i \right] - \sum_{i=1}^{m} \omega_{i} C_{i} \otimes \phi(\mu_{i}) \phi(\mu_{i})^{\top} \right\|_{F}^{2} \\ &= \left\| \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \phi_{T}(t) \phi_{T}(t)^{\top} \middle| G = i \right] - \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \phi_{T}(\mu_{i}) \phi_{T}(\mu_{i})^{\top} \middle| G = i \right] \right\|_{F}^{2} \\ &= \left\| \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \left( \phi_{T}(t) \phi_{T}(t)^{\top} - \phi_{T}(\mu_{i}) \phi_{T}(\mu_{i})^{\top} \right) \middle| G = i \right] \right\|_{F}^{2} \\ &\leq \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \left\| \phi_{X}(x) \phi_{X}(x)^{\top} \otimes \left( \phi_{T}(t) \phi_{T}(t)^{\top} - \phi_{T}(\mu_{i}) \phi_{T}(\mu_{i})^{\top} \right) \right\|_{F}^{2} \middle| G = i \right] \\ &\stackrel{(c)}{=} \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \left\| \phi_{X}(x) \phi_{X}(x)^{\top} \middle| \phi_{T}(t) \phi_{T}(t)^{\top} - \phi_{T}(\mu_{i}) \phi_{T}(\mu_{i})^{\top} \right\|_{F}^{2} \middle| G = i \right] \\ &\stackrel{(c)}{=} \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ \left\| \phi_{T}(t) \phi_{T}(t)^{\top} - \phi_{T}(\mu_{i}) \phi_{T}(\mu_{i})^{\top} \middle| \phi_{F}(\mu_{i})^{\top} \right\|_{F}^{2} \middle| G = i \right] \\ &\stackrel{(c)}{\leq} \sum_{i=1}^{m} \omega_{i} \mathbb{E} \left[ 2 - 2 \exp\left( \frac{-\left\| t - \mu_{i} \right\|_{2}^{2} \middle| G = i \right| \right) \right] \\ &\stackrel{(c)}{\leq} \sum_{i=1}^{m} \omega_{i} \left[ 2 - 2 \exp\left( \frac{-\sigma_{i}^{2}}{\sigma^{2}} \right) \right] \\ &\stackrel{(f)}{\leq} \sum_{i=1}^{m} \omega_{i} \left[ 2 - 2 \exp\left( \frac{-\sigma_{i}^{2}}{\sigma^{2}} \right) \right] \end{aligned}$$

In the above, (a) follows from Jensen's inequality for the convex Frobenius-norm-squared function. (b) holds because  $\|A \otimes B\|_F^2 = \|A\|_F^2 \|B\|_F^2$  for every matrices A, B. (c) comes from the normalized Gaussian kernel satisfying  $\langle \phi_T(t), \phi_T(t) \rangle = k(t,t) = 1$ , resulting in  $\|\phi_T(t)\phi_T(t)^\top\|_F^2 = \text{Tr}(\phi_T(t)\phi_T(t)^\top\phi_T(t)\phi_T(t)^\top) = \text{Tr}(\phi_T(t)\phi_T(t)^\top) = 1$ . (d) follows from the Gaussian kernel definition, proving that  $\phi_T(t)^\top\phi_T(\mu_i) = \exp(-\|t-\mu_i\|_2^2/2\sigma^2)$ . (e) shows the application of Jensen's inequality to the concave  $s(z) = 1 - \exp(-z)$ . (f) holds because  $s(z) = 1 - \exp(-z)$  is a monotonically increasing function. Finally, (g) follows from the inequality  $1 - \exp(-z) \leq z$  for every scalar z. Therefore, the proof is complete.

Next, we apply the Gram–Schmidt process to  $\phi_T(\mu_1), \dots, \phi_T(\mu_m)$  to find orthogonal vectors  $u_1, \dots, u_m$ . We let  $u_1 = \phi_T(\mu_1)$ . Then, for every  $2 \le i \le m$ , we define

$$u_i := \phi(\mu_i) - \sum_{j=1}^{i-1} \langle \phi(\mu_i), u_j \rangle u_j.$$

As a result, the following holds

$$\left\| \sum_{i=1}^{m} \omega_{i} C_{i} \otimes \phi(\mu_{i}) \phi(\mu_{i})^{\top} - \sum_{i=1}^{m} \omega_{i} C_{i} \otimes u_{i} u_{i}^{\top} \right\|_{F}^{2}$$

$$= \left\| \sum_{i=1}^{m} \omega_{i} C_{i} \otimes \left( \phi(\mu_{i}) \phi(\mu_{i})^{\top} - u_{i} u_{i}^{\top} \right) \right\|_{F}^{2}$$

$$\stackrel{(h)}{\leq} \sum_{i=1}^{m} \omega_{i} \left\| C_{i} \otimes \left( \phi(\mu_{i}) \phi(\mu_{i})^{\top} - u_{i} u_{i}^{\top} \right) \right\|_{F}^{2}$$

$$= \sum_{i=1}^{m} \omega_{i} \left\| C_{i} \right\|_{F}^{2} \left\| \phi(\mu_{i}) \phi(\mu_{i})^{\top} - u_{i} u_{i}^{\top} \right\|_{F}^{2}$$

$$\stackrel{(i)}{\leq} \sum_{i=1}^{m} \omega_{i} \left\| \phi(\mu_{i}) \phi(\mu_{i})^{\top} - u_{i} u_{i}^{\top} \right\|_{F}^{2}$$

$$\stackrel{(j)}{=} \sum_{i=1}^{m} \omega_{i} \left( 1 + \left\| u_{i} \right\|^{4} - 2 \left( u_{i}^{\top} \phi_{T}(\mu_{i}) \right)^{2} \right)$$

$$\leq \sum_{i=1}^{m} \omega_{i} \left( 2 - 2 \left( u_{i}^{\top} \phi(\mu_{i}) \right)^{2} \right)$$

$$= 2 \sum_{i=1}^{m} \omega_{i} \left( 1 + u_{i}^{\top} \phi(\mu_{i}) \right) \left( 1 - u_{i}^{\top} \phi(\mu_{i}) \right)$$

$$\stackrel{(k)}{\leq} 4 \sum_{i=1}^{m} \sum_{j=1}^{i-1} \omega_{i} \exp\left( \frac{-\left\| \mu_{i} - \mu_{j} \right\|_{2}^{2}}{\sigma^{2}} \right).$$

Here, (h) follows from the application of Jensen's inequality for the convex Frobenius-norm-squared. (i) holds since the text kernel is normalized and  $\langle \phi_X(x), \phi_X(x) \rangle = k_X(x,x) = 1$ , and therefore  $\|C_i\|_F \leq \mathbb{E}[\|\phi_X(x)\|_2^2] = 1$ . (j) follows from the expansion  $\|uu^\top - vv^\top\|_F^2 = \|u\|_2^4 + \|v\|_2^4 - 2\langle u, v \rangle^2$ . (k) holds because  $u_i^\top \phi_T(\mu_i) \leq 1$  and

$$u_i^{\top} \phi_T(\mu_i) = 1 - \sum_{i=1}^{i-1} \langle \phi_T(\mu_i), u_j \rangle^2 \ge 1 - \sum_{i=1}^{i-1} \exp\left(\frac{-\|\mu_i - \mu_j\|_2^2}{\sigma^2}\right).$$

Since we know that for every matrices  $A, B \in \mathbb{R}^{d \times d}$ ,  $\|A + B\|_F^2 \le 2\|A\|_F^2 + 2\|B\|_F^2$ , the above results show that

$$\left\| C_{X \otimes T} - \sum_{i=1}^{m} \omega_i C_i \otimes u_i u_i^{\top} \right\|_F^2 \le \sum_{i=1}^{m} 4\omega_i \frac{\sigma_i^2}{\sigma^2} + \sum_{i=2}^{m} \sum_{j=1}^{i-1} 8\omega_i \exp\left(\frac{-\|\mu_i - \mu_j\|_2^2}{\sigma^2}\right).$$

As a result, the Hoffman-Wielandt inequality shows that for the sorted eigenvalues vector  $\widetilde{\boldsymbol{\lambda}}$  of  $C_{X\otimes T}$  and sorted eigenvalues vector  $\widetilde{\boldsymbol{\lambda}}$  of  $\sum_{i=1}^m \omega_i C_i \otimes u_i u_i^{\top}$  the following holds:

$$\|\boldsymbol{\lambda} - \widetilde{\boldsymbol{\lambda}}\|_{2}^{2} \leq \|C_{X \otimes T} - \sum_{i=1}^{m} \omega_{i} C_{i} \otimes u_{i} u_{i}^{\top}\|_{F}^{2}$$

$$\leq \sum_{i=1}^{m} 4\omega_{i} \frac{\sigma_{i}^{2}}{\sigma^{2}} + \sum_{i=2}^{m} \sum_{j=1}^{i-1} 8\omega_{i} \exp\left(\frac{-\|\mu_{i} - \mu_{j}\|_{2}^{2}}{\sigma^{2}}\right).$$

Since  $u_1,\ldots,u_m$  are orthogonal vectors, the definition of Kronecker product implies that the eigenvalues of  $\sum_{i=1}^m \omega_i C_i \otimes u_i u_i^{\top}$  will be the union of the eigenvalues of  $\omega_i C_i \otimes u_i u_i^{\top}$  over  $i \in \{1,\ldots,m\}$ . On the other hand, we know that the non-zero eigenvalues of  $\omega_i C_i \otimes u_i u_i^{\top}$  will be equal to the factor  $\omega_i \|u_i\|_2^2$  times the eigenvalues of  $C_i$ . Also, we know that  $1 \geq \|u_i\|_2^2 \geq 1 - 2\sum_{j=1}^{i-1} \exp(-\frac{\|\mu_i - \mu_j\|_2^2}{\sigma^2})$ . Consequently, we can show that for vector  $\widehat{\lambda}_{x \otimes t} = \operatorname{Union}(\omega_i \operatorname{Eigs}(C_i): i \in \{1,\ldots,m\})$ , we have the following for every  $\alpha \geq 2$  and defined increasing function g in Theorem 1

$$\begin{split} \left| \widetilde{H}_{\alpha}(X,T) - \frac{1}{1-\alpha} \log(\|\widehat{\lambda}_{x \otimes t}\|_{\alpha}^{\alpha}) \right| &\leq g(\|\widetilde{\lambda}_{x \otimes t}\|_{\alpha} - \|\widehat{\lambda}_{x \otimes t}\|_{\alpha}) \\ &\leq g(\|\operatorname{sort}(\widetilde{\lambda}_{x \otimes t}) - \operatorname{sort}(\widehat{\lambda}_{x \otimes t})\|_{\alpha}) \\ &\leq g(\|\operatorname{sort}(\widetilde{\lambda}_{x \otimes t}) - \operatorname{sort}(\widehat{\lambda}_{x \otimes t})\|_{2}) \\ &\leq g(\sum_{i=1}^{m} 4\omega_{i} \frac{\sigma_{i}^{2}}{\sigma^{2}} + \sum_{i=2}^{m} \sum_{j=1}^{i-1} 16\omega_{i} \exp\left(\frac{-\|\mu_{i} - \mu_{j}\|_{2}^{2}}{\sigma^{2}}\right)). \end{split}$$

Note that the above proof holds for every marginal distribution on X, and we choose a deterministic constant  $X = \mathbf{0}$ , then the joint entropy reduces to the marginal entropy and the above inequality also shows the following:

$$\left| \widetilde{H}_{\alpha}(T) - \frac{1}{1-\alpha} \log \left( \| [\omega_1, \dots, \omega_m] \|_{\alpha}^{\alpha} \right) \right| \leq g \left( \sum_{i=1}^m 4\omega_i \frac{\sigma_i^2}{\sigma^2} + \sum_{i=2}^m \sum_{j=1}^{i-1} 16\omega_i \exp \left( \frac{-\|\mu_i - \mu_j\|_2^2}{\sigma^2} \right) \right).$$

Therefore, following the Triangle inequality and the definition  $\widetilde{H}_{\alpha}(X|T) = \widetilde{H}_{\alpha}(X,T) - \widetilde{H}_{\alpha}(T)$ , the previous two inequalities prove that

$$\left| \widetilde{H}_{\alpha}(X|T) - \left( \frac{1}{1-\alpha} \log(\|\widehat{\lambda}_{x \otimes t}\|_{\alpha}^{\alpha}) - \frac{1}{1-\alpha} \log(\|[\omega_{1}, \dots, \omega_{m}]\|_{\alpha}^{\alpha}) \right) \right|$$

$$\leq 2g \left( \sum_{i=1}^{m} 4\omega_{i} \frac{\sigma_{i}^{2}}{\sigma^{2}} + \sum_{i=2}^{m} \sum_{j=1}^{i-1} 16\omega_{i} \exp\left( \frac{-\|\mu_{i} - \mu_{j}\|_{2}^{2}}{\sigma^{2}} \right) \right).$$

On the other hand, we can simplify the above expression as

$$\frac{1}{1-\alpha} \log(\|\widehat{\lambda}_{x \otimes t}\|_{\alpha}^{\alpha}) - \frac{1}{1-\alpha} \log(\|[\omega_{1}, \dots, \omega_{m}]\|_{\alpha}^{\alpha})$$

$$= \frac{1}{1-\alpha} \log(\sum_{i=1}^{m} \omega_{i}^{\alpha} \|\lambda_{C_{i}}\|_{\alpha}^{\alpha}) - \frac{1}{1-\alpha} \log(\sum_{i=1}^{m} \omega_{i}^{\alpha})$$

$$= \frac{1}{1-\alpha} \log(\sum_{i=1}^{m} \frac{\omega_{i}^{\alpha}}{\sum_{j=1}^{m} \omega_{j}^{\alpha}} \|\lambda_{C_{i}}\|_{\alpha}^{\alpha})$$

Note that the definition  $f_{\alpha}(t) = \exp((1-\alpha)t)$  implies that  $f_{\alpha}^{-1}(z) = \frac{1}{1-\alpha}\log(z)$ , which connects to the entropy definition as  $H(X|G=i) = f_{\alpha}^{-1}(\|\lambda_{C_i}\|_{\alpha}^{\alpha})$ . As a result, we can combine the previous two equations and complete the proof as:

$$\left| \widetilde{H}_{\alpha}(X|T) - f_{\alpha}^{-1} \left( \sum_{i=1}^{m} \frac{\omega_{i}^{\alpha}}{\sum_{j=1}^{m} \omega_{j}^{\alpha}} f_{\alpha} \left( \widetilde{H}_{\alpha}(X|G=i) \right) \right) \right|$$

$$\leq 2g \left( \sum_{i=1}^{m} 4\omega_{i} \frac{\sigma_{i}^{2}}{\sigma^{2}} + \sum_{i=2}^{m} \sum_{j=1}^{i-1} 16\omega_{i} \exp\left( \frac{-\|\mu_{i} - \mu_{j}\|_{2}^{2}}{\sigma^{2}} \right) \right).$$

#### **D. Additional Numerical Results**

#### D.1. Measuring Conditional-Vendi across prompt types

In this section, we conducted additional experiments similar to those in Figure 4. We created 5,000 prompts across different categories using GPT40 and generated corresponding images with text-to-image models. We reported Conditional-Vendi for the top 3 groups in the text data on PixArt- $\alpha$ , Stable Diffusion XL text-to-image generative models.

As shown in Figure 7, Figure 8 and Figure 9, we observed the same behavior during these experiments: the Conditional-Vendi score for "dog" prompts was significantly higher than for the "airplane" and "sofa" categories. This observation suggests that the outputs of generative models are unbalanced when presented with different groups of text prompts.

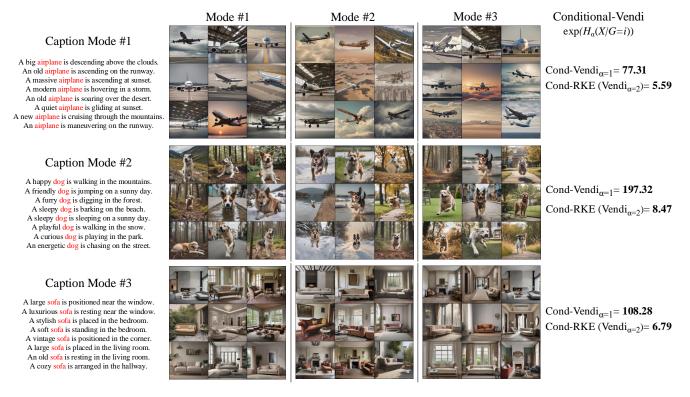


Figure 7. Quantifying image diversity for different clusters of text prompts. Images are generated using the Stable Diffusion XL model.

#### D.2. Quantifying model-induced diversity via Conditional-Vendi.

In this section, we provided a more detailed version of Figure 2. As shown in Figure 10, we found that Conditional-Vendi increased at a more rapid rate when the prompts did not specify the type of animal in the picture. In contrast, when the animal types were specified in the prompts, there was only a slight increase in the Conditional-Vendi score.

#### D.3. Additional Numerical Evaluation of the Conditional-Vendi Score

To further experiment the correlation between the intrinsic model diversity and the defined Conditional-Vendi score, we have performed experiments of quantifying the diversity scores for unspecified and type-specified prompts when generating data from standard text-to-image models. To conduct an extensive evaluation of the Conditional-Vendi score, we performed the experiments on the nine combinations of three category types, 1) animals, 2) fruits, 3) objects, and three SOTA text-to-image generation models SDXL, Kandinsky, and PixArt- $\Sigma$  Chen et al. [9]. In each of the nine experiments, we generated prompts on 10 different types related to each category and created image samples by inputting the prompts to the text-to-image model. In each experiment, we simulated 10 prompt-based generative models by considering image samples from j types for  $j \in \{1, \ldots, 10\}$ .

To validate the correlation between the Conditional-Vendi score and the groundtruth non-prompt-induced diversity, we evaluated the conditional-Vendi scores of the generated samples when the type is specified in the prompt (i.e. the original prompt) and when the type is unspecified by the prompt. Figures 13,15,14 show the Conditional-Vendi scores evaluated under type-specified and unspecified. In these figures, we validate that by specifying the category type in the prompt, the conditional Vendi score will have a lower value and grows slowly with the number of types expressed in the prompts. In addition, given the original prompts specifying the type of category in the image, we compared the Vendi and Conditional-Vendi scores among the three text-to-image models. Figures ??? show the comparison between the scores of the three models, which suggest the higher intrinsic diversity measure by Conditional-Vendi for the SD-XL and PixArt- $\Sigma$  models.

#### Mode #1 Mode #2 Mode #3 Conditional-Vendi $\exp(H_{\alpha}(X/G=i))$ Caption Mode #1 A big aimplane is descending above the clouds. Cond-Vendi<sub> $\alpha=1$ </sub>= **21.29** An old aimlane is ascending on the runway. A massive aimplane is ascending at sunset. Cond-RKE (Vendi<sub> $\alpha=2$ </sub>)= **3.07** A modern airplane is hovering in a storm. An old airplane is soaring over the desert. A quiet aimplane is gliding at sunset. A new aimplane is cruising through the mountains An airplane is maneuvering on the runway Caption Mode #2 A happy dog is walking in the mountains. Cond-Vendi<sub> $\alpha=1$ </sub>= **82.81** A friendly dog is jumping on a sunny day. A furry dog is digging in the forest. A sleepy dog is barking on the beach. A sleepy dog is sleeping on a sunny day. Cond-RKE (Vendi<sub> $\alpha=2$ </sub>)= **6.35** A playful dog is walking in the snow A curious dog is playing in the park. An energetic dog is chasing on the street. Caption Mode #3 A large sofa is positioned near the window. A luxurious sofa is resting near the window. A stylish sofa is placed in the bedroom. Cond-Vendi<sub> $\alpha=1$ </sub>= **36.85** Cond-RKE (Vendi<sub> $\alpha=2$ </sub>)= **4.16** A soft sofa is standing in the bedroom. A vintage sofa is positioned in the corner. A large sofa is placed in the living room. An old sofa is resting in the living room. A cozy sofa is arranged in the hallway.

Figure 8. Quantifying image diversity for different clusters of text prompts. Images are generated using the PixArt- $\alpha$  model.

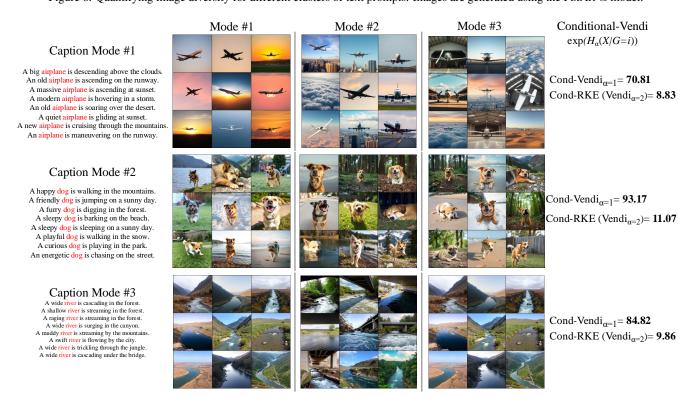


Figure 9. Quantifying image diversity for different clusters of text prompts. Images are generated using the Flux model.

### Type-specified animal prompts

### Prompts:

A fox is rolling in the grass.

A camel is resting near a sand dune.

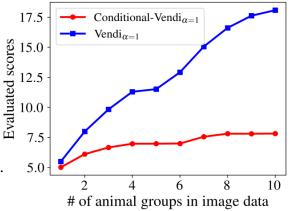
A wolf is drinking from a river.

A cow is resting under a tree.

A sheep is resting in a grassy pasture.

An elephant is walking through thick jungle.

A giraffe is reaching up for leaves.



# Unspecified animal prompts

### Prompts:

An animal is rolling in the grass.

An animal resting near a sand dune.

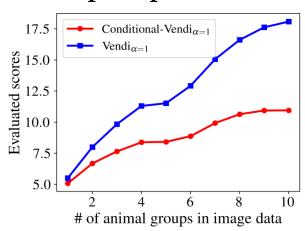
An animal is drinking from a river.

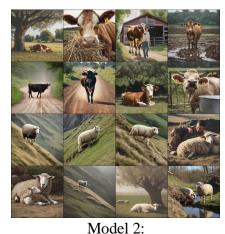
An animal is resting under a tree.

An animal is resting in a grassy pasture.

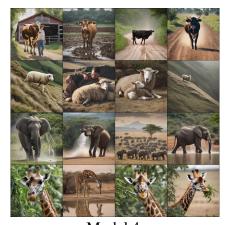
An animal is walking through jungle.

An animal is reaching up for leaves.





Samples from 2 animal groups



Model 4: Samples from 4 animal groups



Model 8: Samples from 8 animal groups

Figure 10. Comparing Conditional-Vendi with Vendi on different animal groups generated by Stable Diffusion-XL.

## Type-specified animal prompts

### Prompts:

A fox is rolling in the grass.

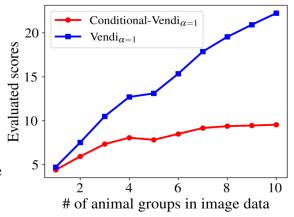
A camel is resting near a sand dune.

A wolf is drinking from a river.

A cow is resting under a tree.

A sheep is resting in a grassy pasture.

An elephant is walking through thick jungle A giraffe is reaching up for leaves.



### Unspecified animal prompts

### Prompts:

An animal is rolling in the grass.

An animal resting near a sand dune.

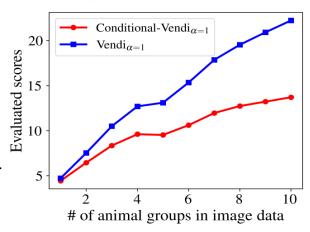
An animal is drinking from a river.

An animal is resting under a tree.

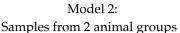
An animal is resting in a grassy pasture.

An animal is walking through jungle.

An animal is reaching up for leaves.









Model 4: Samples from 4 animal groups



Model 8: Samples from 8 animal groups

Figure 11. Comparing Conditional-Vendi with Vendi on different animal groups generated by  $PixArt\Sigma$ .

## Type-specified animal prompts

### Prompts:

A fox is rolling in the grass.

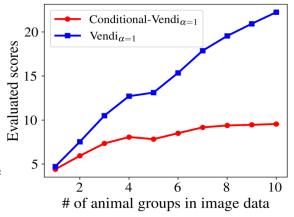
A camel is resting near a sand dune.

A wolf is drinking from a river.

A cow is resting under a tree.

A sheep is resting in a grassy pasture.

An elephant is walking through thick jungle A giraffe is reaching up for leaves.



### Unspecified animal prompts

### Prompts:

An animal is rolling in the grass.

An animal resting near a sand dune.

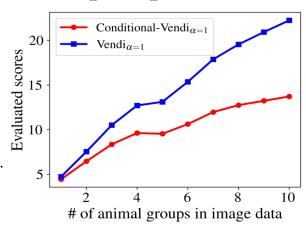
An animal is drinking from a river.

An animal is resting under a tree.

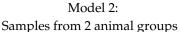
An animal is resting in a grassy pasture.

An animal is walking through jungle.

An animal is reaching up for leaves.









Model 4: Samples from 4 animal groups



Model 8: Samples from 8 animal groups

Figure 12. Comparing Conditional-Vendi with Vendi on different animal groups generated by Kandinsky.

### Prompts:

A Chair is placed under a sprawling tree.

A Sofa is glowing in the light of a nearby fire.

A Book is balancing on the edge of a table.

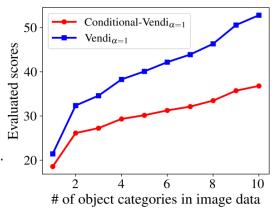
A Clock is positioned in an office setup.

A Lamp is sitting under a hanging light bulb.

A laptop is half-hidden behind a stack of boxes.

A car is in the corner of a large warehouse.

The cup is precariously balanced on rocks.



### Unspecified object prompts

### Prompts:

An **object** is placed under a sprawling tree.

An **object** is glowing in the light of a nearby fire.

An **object** is balancing on the edge of a table.

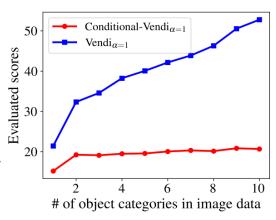
An **object** is positioned in an office setup.

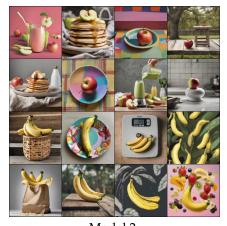
An **object** is sitting under a hanging light bulb.

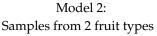
An **object** is half-hidden behind a stack of boxes.

An **object** is in the corner of a large warehouse.

The **object** is precariously balanced on rocks.









Model 4: Samples from 4 fruit types



Model 8: Samples from 8 fruit types

Figure 13. Comparing Conditional-Vendi with Vendi on different fruit types generated by Stable Diffusion-XL.

#### D.4. Correlation between GroundTruth-Cluster-Vendi and Conditional-Vendi Scores

To validate the theoretical connection between the Vendi and Conditional-Vendi scores, we performed an experiment and evaluated a baseline metric called GroundTruth-Cluster-Vendi score. To measure the GroundTruth-Cluster-Vendi score, we

utilize the side knowledge of the ground-truth clusters of the input prompts and then compute and average the regular Vendi scores for the data generated within each cluster. Mathematically, given t sample cluster sets in  $S = \{S_1, \ldots, S_t\}$ , which partition the input text indices  $\{1, \ldots, n\}$ , we define the Cluster-Vendi score as follows, where  $|S_j|$  denotes the cardinality of subset  $S_j$ :

Cluster-Vendi
$$(x_1, \ldots, x_n \mid S) := \sum_{i=1}^t \frac{|S_i|}{n} \cdot \text{Vendi}(\{x_j : j \in S_i\}).$$

Note that the above definition requires the knowledge of the clusters, which could be given by an oracle in the case of the GroundTruth-Cluster-Vendi score, or computed by a clustering algorithm such as K-Means to obtain the KMeans-Cluster-Vendi score. Observe that given the knowledge of the clusters revealed by an oracle, the GroundTruth-Cluster-Vendi score is a sensible definition of internal model diversity, which, as shown in Theorem 1, is expected to correlate with our defined Conditional-Vendi score.

In the numerical settings of the previous section, where we know the ground-truth clusters based on the type of animal, fruit, or object in the texts, we computed the value of the GroundTruth-Cluster-Vendi score and compared it with the evaluated Conditional-Vendi score. As demonstrated in Figures ??, the two diversity scores, Conditional-Vendi and GroundTruth-Cluster-Vendi, highly correlate for the ten simulated generative models in the experiments.

However, note that in a real-world scenario, we do not have access to the ground-truth clusters. To estimate the score, we should use a clustering algorithm such as K-Means to find the clusters and compute the Cluster-Vendi score. We note that the optimization problem addressed by standard clustering algorithms represents a challenging non-convex optimization, which, depending on the algorithm's initial point, could converge to different solutions. Our numerical results with the K-Means clustering algorithm in Figure 19 also demonstrated these clustering challenges and, in several cases, failed to find the ground-truth clusters with high accuracy.

#### D.5. Algorithm for Computing Conditional-Vendi and Information-Vendi

In this section, we present the algorithm to compute the Conditional-Vendi and Information-Vendi scores. Using the definition provided in Section 4, combined with the entropy definition in (3), we calculate the Conditional-Vendi score. The steps are outlined in Algorithm 1.

#### D.6. Qualitative results for generative models trained on MS-COCO dataset

In this section, we provide images and prompts corresponding to Figure 3. Figure 21 illustrates three clusters obtained by applying KMeans to cluster MS-COCO validation set prompts into 1000 clusters. The images are presented for four generative models. Comparing the prompts with the generated images reveals that FLUX exhibits the highest alignment between text and image, while GigaGAN demonstrates greater diversity but misses some features of the prompts. These observations are further supported by the Conditional-Vendi and Information-Vendi metrics.

#### D.7. Effect of Bandwidth on Conditional-Vendi and Information-Vendi

To further investigate the effect of bandwidth on Conditional-Vendi and Information-Vendi, we began by selecting the image bandwidth similar to prior works Friedman and Dieng [11], Ospanov et al. [31]. We then measured and plotted the scores using varying text kernel bandwidths. Figure 22 demonstrates consistent rankings of the four models across different bandwidth parameters. The results indicate that as the kernel bandwidth increases, the number of text clusters increases, leading to a decrease in the Information-Vendi value.

#### Algorithm 1 Conditional-Vendi and Information-Vendi

- 1: Input: Sample sets  $\{\mathbf{x}_1,\dots,\mathbf{x}_n\}$  and  $\{\mathbf{t}_1,\dots,\mathbf{t}_n\}$ , Gaussian kernel bandwidths  $\sigma_i^2,\sigma_t^2$ , order  $\alpha$ . 2: Compute kernel matrices:  $K_{\mathbf{X}} = \frac{1}{n}[k(\mathbf{x}_i,\mathbf{x}_j)]_{n\times n}, K_{\mathbf{T}} = \frac{1}{n}[k(\mathbf{t}_i,\mathbf{t}_j)]_{n\times n}$ 3: Perform eigendecomposition on the  $K_{\mathbf{X}},K_{\mathbf{T}}$  and  $\frac{1}{n}K_{\mathbf{X}}\odot K_{\mathbf{T}}$  matrices:

$$\begin{aligned} \{\lambda_1^{\mathbf{X}}, \dots, \lambda_n^{\mathbf{X}}\} &\leftarrow \operatorname{Eigendecomposition}(K_{\mathbf{X}}) \\ \{\lambda_1^{\mathbf{T}}, \dots, \lambda_n^{\mathbf{T}}\} &\leftarrow \operatorname{Eigendecomposition}(K_{\mathbf{T}}) \\ \{\lambda_1^{\mathbf{X}, \mathbf{T}}, \dots, \lambda_n^{\mathbf{X}, \mathbf{T}}\} &\leftarrow \operatorname{Eigendecomposition}(\frac{1}{n}K_{\mathbf{X}} \odot K_{\mathbf{T}}) \end{aligned}$$

4: Compute  $H_{\alpha}(\frac{1}{n}K_{\mathbf{X}})$ ,  $H_{\alpha}(\frac{1}{n}K_{\mathbf{T}})$  and  $H_{\alpha}(\frac{1}{n}K_{\mathbf{X}}\odot K_{\mathbf{T}})$  using their eigenvalues.

$$H_{\alpha}\left(\frac{1}{n}K_{\mathbf{X}}\right) \leftarrow \frac{1}{1-\alpha}\log\left(\sum_{i=1}^{n}(\lambda_{i}^{\mathbf{X}})^{\alpha}\right)$$

$$H_{\alpha}\left(\frac{1}{n}K_{\mathbf{T}}\right) \leftarrow \frac{1}{1-\alpha}\log\left(\sum_{i=1}^{n}(\lambda_{i}^{\mathbf{T}})^{\alpha}\right)$$

$$H_{\alpha}\left(\frac{1}{n}K_{\mathbf{X}}\odot K_{\mathbf{T}}\right) \leftarrow \frac{1}{1-\alpha}\log\left(\sum_{i=1}^{n}(\lambda_{i}^{\mathbf{X},\mathbf{T}})^{\alpha}\right)$$

5: Compute Conditional-Vendi and Information-Vendi

Conditional-Vendi
$$_{\alpha}(x_1, \dots, x_n | t_1, \dots, t_n) \leftarrow \exp\left(H_{\alpha}\left(\frac{1}{n}K_{\mathbf{X}} \odot K_{\mathbf{T}}\right) - H_{\alpha}\left(\frac{1}{n}K_{\mathbf{T}}\right)\right)$$
  
Information-Vendi $_{\alpha}(x_1, \dots, x_n; t_1, \dots, t_n) \leftarrow \exp\left(H_{\alpha}\left(\frac{1}{n}K_{\mathbf{X}}\right) + H_{\alpha}\left(\frac{1}{n}K_{\mathbf{T}}\right) - H_{\alpha}\left(\frac{1}{n}K_{\mathbf{X}} \odot K_{\mathbf{T}}\right)\right)$ 

6: Output: Conditional-Vendi and Information-Vendi

### Type-specified fruit prompts

Prompts:

An **apple** is next to a cold glass of fresh juice.

A **banana** is being sliced with a sharp knife.

The **watermelon** is blended into a smoothie.

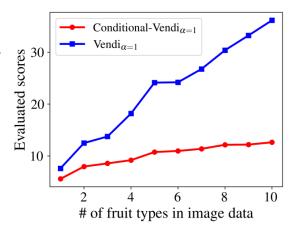
The **pineapple** is falling out of a grocery bag.

A **strawberry** is being washed.

The **peach** is sitting on a kitchen countertop.

A **cherry** is being sliced with a knife.

A **mango** is sitting on a colorful plate.



### Unspecified fruit prompts

Prompts:

A fruit is next to a cold glass of fresh juice.

A fruit is being sliced with a sharp knife.

The fruit is blended into a smoothie.

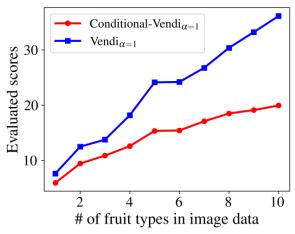
The fruit is falling out of a grocery bag.

A fruit is being washed.

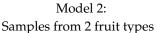
The **fruit** is sitting on a kitchen countertop.

A **fruit** is being sliced with a knife.

A **fruit** is sitting on a colorful plate.









Model 4: Samples from 4 fruit types



Model 8: Samples from 8 fruit types

Figure 14. Comparing Conditional-Vendi with Vendi on different fruit types generated by Kandinsky.

### Type-specified fruit prompts

Prompts:

An **apple** is next to a cold glass of fresh juice.

A **banana** is being sliced with a sharp knife.

The **watermelon** is blended into a smoothie.

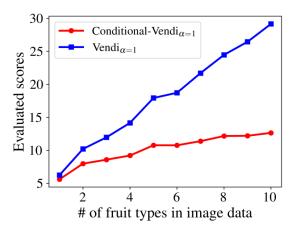
The **pineapple** is falling out of a grocery bag.

A **strawberry** is being washed.

The **peach** is sitting on a kitchen countertop.

A **cherry** is being sliced with a knife.

A **mango** is sitting on a colorful plate.



### Unspecified fruit prompts

Prompts:

A **fruit** is next to a cold glass of fresh juice.

A **fruit** is being sliced with a sharp knife.

The **fruit** is blended into a smoothie.

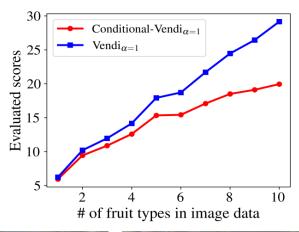
The **fruit** is falling out of a grocery bag.

A **fruit** is being washed.

The **fruit** is sitting on a kitchen countertop.

A **fruit** is being sliced with a knife.

A **fruit** is sitting on a colorful plate.





Model 2: Samples from 2 fruit types



Model 4: Samples from 4 fruit types



Model 8: Samples from 8 fruit types

Figure 15. Comparing Conditional-Vendi with Vendi on different fruit types generated by PixArt- $\Sigma$ .

### Prompts:

A Chair is placed under a sprawling tree.

A Sofa is glowing in the light of a nearby fire.

A Book is balancing on the edge of a table.

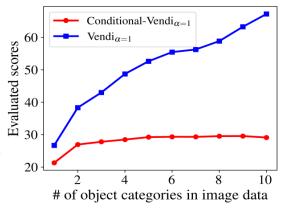
A Clock is positioned in an office setup.

A Lamp is sitting under a hanging light bulb.

A laptop is half-hidden behind a stack of boxes.

A car is in the corner of a large warehouse.

The cup is precariously balanced on rocks.



### Unspecified object prompts

### Prompts:

An **object** is placed under a sprawling tree.

An **object** is glowing in the light of a nearby fire.

An **object** is balancing on the edge of a table.

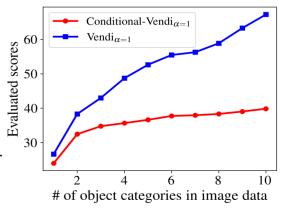
An **object** is positioned in an office setup.

An **object** is sitting under a hanging light bulb.

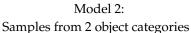
An **object** is half-hidden behind a stack of boxes.

An **object** is in the corner of a large warehouse.

The **object** is precariously balanced on rocks.









Model 4: Samples from 4 object categories



Model 8: Samples from 8 object categories

Figure 16. Comparing Conditional-Vendi with Vendi on different fruit types generated by Stable Diffusion-XL.

## Prompts:

A **Chair** is placed under a sprawling tree.

A **Sofa** is glowing in the light of a nearby fire.

A **Book** is balancing on the edge of a table.

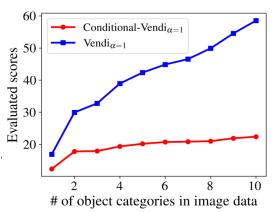
A **Clock** is positioned in an office setup.

A **Lamp** is sitting under a hanging light bulb.

A **laptop** is half-hidden behind a stack of boxes.

A **car** is in the corner of a large warehouse.

The **cup** is precariously balanced on rocks.



### Unspecified object prompts

### Prompts:

An **object** is placed under a sprawling tree.

An **object** is glowing in the light of a nearby fire.

An **object** is balancing on the edge of a table.

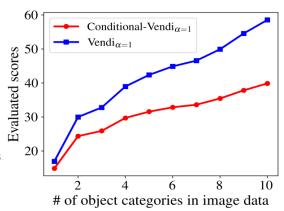
An **object** is positioned in an office setup.

An **object** is sitting under a hanging light bulb.

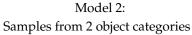
An **object** is half-hidden behind a stack of boxes

An **object** is in the corner of a large warehouse.

The **object** is precariously balanced on rocks.









Model 4: Samples from 4 object categories



Model 8: Samples from 8 object categories

Figure 17. Comparing Conditional-Vendi with Vendi on different object categories types generated by Kandinsky.

### Prompts:

A **Chair** is placed under a sprawling tree.

A **Sofa** is glowing in the light of a nearby fire.

A **Book** is balancing on the edge of a table.

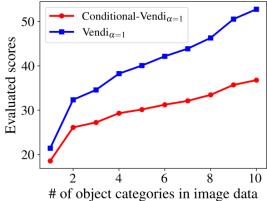
A **Clock** is positioned in an office setup.

A **Lamp** is sitting under a hanging light bulb.

A **laptop** is half-hidden behind a stack of boxes.

A **car** is in the corner of a large warehouse.

The **cup** is precariously balanced on rocks.



### Unspecified object prompts

### Prompts:

An **object** is placed under a sprawling tree.

An **object** is glowing in the light of a nearby fire.

An **object** is balancing on the edge of a table.

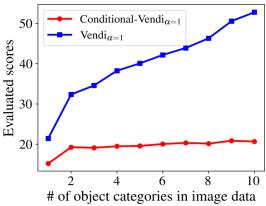
An **object** is positioned in an office setup.

An **object** is sitting under a hanging light bulb.

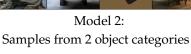
An **object** is half-hidden behind a stack of boxes.

An **object** is in the corner of a large warehouse.

The **object** is precariously balanced on rocks.









Model 4: Samples from 4 object categories



Model 8: Samples from 8 object categories

Figure 18. Comparing Conditional-Vendi with Vendi on different object categories types generated by  $PixArt-\Sigma$ .

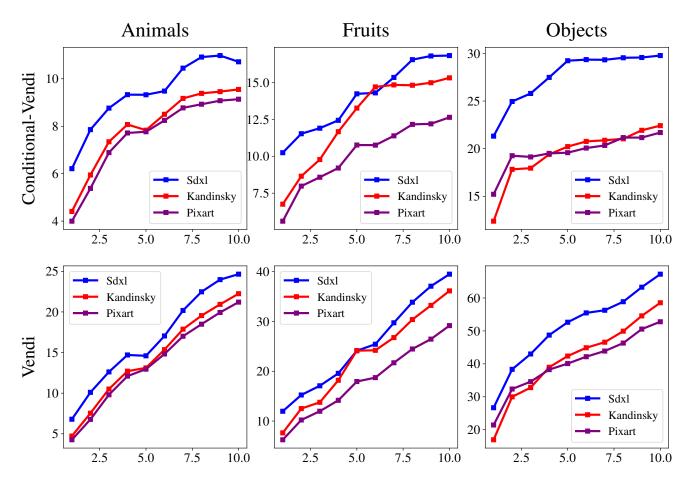
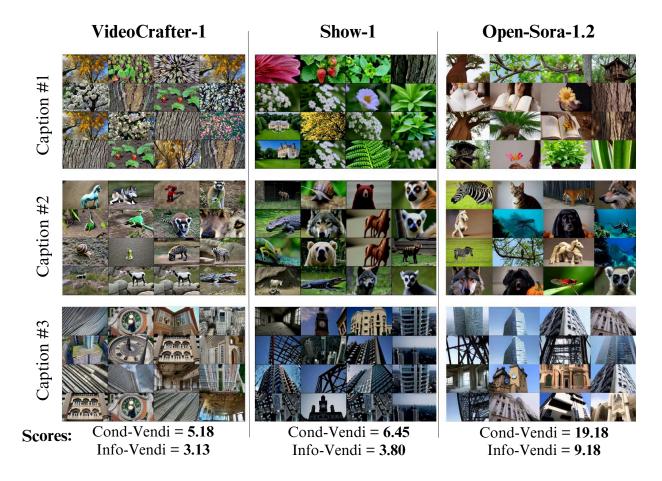


Figure 19. Evaluated (unconditional) Vendi and Conditional-Vendi scores of three text-to-image models in the category-based experiments with varying number of types within each of the categories: Animals, Fruits, Objects.



### Caption Mode #1

close up video of flower petals curious cat sitting and looking around high angle shot of a clock tower a leaf on a glass the long trunks of tall trees in the forest trees in the forest during sunny day close up video of tree bark reflection of tree branches trunks of many trees in the forest tree leaves providing shades from the sun leaves swaying in the wind low angle shot of baobab tree close up video of strawberry plant close up video of tree bark tree with golden leaves close up view of a plant

### Caption Mode #2

# Caption Mode #3

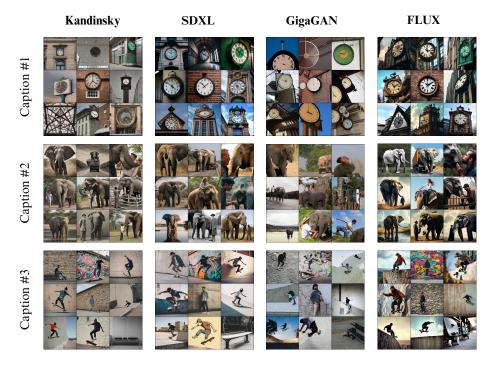
close up shot of a wild bear a zebra eating grass on the field a bear wearing red jersey close up video of snail a gorilla eating a carrot close up of wolf a meerkat looking around a hyena in a zoo lemur eating grass leaves an owl being trained by a man an American crocodile

a black dog wearing halloween costume close up shot of a steel structure an apartment building with balcony low angle shot of a building tower on hill a miniature house eiffel tower from the seine river low angle footage of an apartment building

island with pier and antique building asian historic architecture mosque in the middle east low angle shot of a building top view of a high rise building

Figure 20. Measuring Conditional-Vendi and Information-Vendi for text-to-video models

close up of a lemur



#### Caption Mode #1

A building displaying a clock showing the time to be 6 oclock.

A clock hanging from the ceiling of a building.

A large metal green clock hanging from the side of a building.

A clock that is on top of a sign.

A large clock mounted to a brick wall.

A large clock hanging off the side of a tall building.

A clock in near the triangular roof of a large building.

A large clock and a sign on top of a building.

A large clock mounted to the side of a building.

#### Caption Mode #2

The elephant has a large white spot on its abdomen.

The truck driver hauls an elephant down the highway.

A man getting a kiss on the neck from an elephant's trunk

A large elephant walking next to a man

A woman in white shirt climbing onto an elephant.

A man is leaning over a fence offering food to an elephant/

A large elephant standing on the side of a lake.

A man standing next to an elephant who stole his hat with it's trunk.

A man standing near an elephant with its trunk outstretched.

#### Caption Mode #3

A young man riding a skateboard on a stone wall.

A man balancing on a skateboard in front of a graffiti covered wall.

A man doing a trick on a wall with a skateboard.

Bearded skateboarder maintains balance while skating up wall.

A man standing next to a stone wall while holding a skateboard.

There is a man skateboard on the side of a wall.

a guy skate boarding on the edge of a wall

A man on a skateboard is trying to jump over a wall.

two black and white skate boards under a black steel bench

Figure 21. Effect of text kernel bandwidth on Conditional-Vendi and Information-Vendi scores

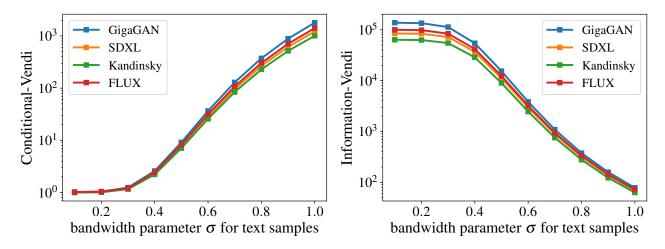


Figure 22. Effect of text kernel bandwidth on Conditional-Vendi and Information-Vendi scores