

OTPrune: Distribution-Aligned Visual Token Pruning via Optimal Transport (Supplementary Material)

Xiwen Chen^{1,2*} Wenhui Zhu^{3*†} Gen Li² Xuanzhao Dong³ Yujian Xiong³ Hao Wang²
 Peijie Qiu⁴ Qingquan Song⁵ Zhipeng Wang^{6†‡} Shao Tang⁷ Yalin Wang³ Abolfazl Razi^{2§}
¹ Morgan Stanley, ² Clemson University, ³ Arizona State University,
⁴ Washington University in St. Louis, ⁵ Texas A&M University,
⁶ Rice University ⁷ Florida State University

A. Proofs

A.1. Proof of Inequality

Recap, for a positive semidefinite square matrix \mathbf{X} , the operator $\Psi(\cdot)$ is presented as:

$$\Psi(\mathbf{X}) = \log \det(\mathbf{I} + \gamma \mathbf{X}) = \sum_{i=1}^d \log(1 + \gamma \lambda_i), \quad (1)$$

where $\{\lambda_i\}$ are the eigenvalues of \mathbf{X} and $\gamma > 0$ is a fixed scaling parameter. Since \mathbf{X} is positive semidefinite (so $\lambda_i \geq 0$), we have the following inequalities,

$$\Psi(\mathbf{X}) \leq \gamma \operatorname{tr}(\mathbf{X}) \leq \gamma (\operatorname{tr}(\mathbf{X}^{1/2}))^2. \quad (2)$$

Now, we will show its proof as follows.

Proof. Let $\mathbf{X} \succeq \mathbf{0}$ with non-negative eigenvalues $\{\lambda_i\}$. For the first inequality in (2), use the elementary bound $\log(1 + t) \leq t$ for all $t \geq 0$. Since $\gamma \lambda_i \geq 0$,

$$\log(1 + \gamma \lambda_i) \leq \gamma \lambda_i \quad \Rightarrow \quad \Psi(\mathbf{X}) \leq \gamma \sum_{i=1}^d \lambda_i = \gamma \operatorname{tr}(\mathbf{X}).$$

For the second inequality in (2), note that

$$\begin{aligned} (\operatorname{tr}(\mathbf{X}^{1/2}))^2 &= \left(\sum_{i=1}^d \sqrt{\lambda_i} \right)^2 \\ &= \sum_{i=1}^d \lambda_i + 2 \sum_{1 \leq i < j \leq d} \sqrt{\lambda_i \lambda_j} \\ &\geq \sum_{i=1}^d \lambda_i = \operatorname{tr}(\mathbf{X}), \end{aligned}$$

because each term $\sqrt{\lambda_i \lambda_j}$ is nonnegative. Multiplying both sides by $\gamma > 0$ gives

$$\gamma \operatorname{tr}(\mathbf{X}) \leq \gamma (\operatorname{tr}(\mathbf{X}^{1/2}))^2.$$

Combining the two inequalities yields

$$\Psi(\mathbf{X}) \leq \gamma \operatorname{tr}(\mathbf{X}) \leq \gamma (\operatorname{tr}(\mathbf{X}^{1/2}))^2,$$

as claimed. □

A.2. Proof of Submodularity

Definition 1 (Submodularity). A set function $h : 2^{[m]} \rightarrow \mathbb{R}$ is *submodular* if it satisfies the diminishing-returns property: for any $S \subseteq T \subseteq [m]$ and any $i \notin T$,

$$h(S \cup \{i\}) - h(S) \geq h(T \cup \{i\}) - h(T).$$

Lemma 1. ([1–3]) Let $\mathbf{A} \in \mathbb{R}^{m \times m}$ be a positive definite matrix, and for any subset $S \subseteq \{1, \dots, m\}$ let \mathbf{A}_S denote the corresponding principal submatrix. Define the set function

$$h(S) = \log \det(\mathbf{A}_S).$$

Then h is a submodular set function.

Proposition 2. Let $\mathbf{V} \in \mathbb{R}^{m \times d}$ and let $\mathbf{V}_{\mathcal{C}}$ denote the submatrix of \mathbf{V} consisting of rows indexed by $\mathcal{C} \subseteq [m]$. Define

$$\Gamma(\mathcal{C}) = \log \det(\mathbf{I} + \tilde{\gamma} \mathbf{V}_{\mathcal{C}} \mathbf{V}^{\top} \mathbf{V} \mathbf{V}_{\mathcal{C}}^{\top}), \quad \tilde{\gamma} > 0.$$

Then Γ is normalized and submodular.

Proof. To prove this proposition, We need to first prove this: Let $\mathbf{V} \in \mathbb{R}^{m \times d}$ be any matrix and $\mathcal{C} \subseteq [m]$ be an index set. Let $\mathbf{V}_{\mathcal{C}}$ denote the submatrix of \mathbf{V} formed by the rows indexed by \mathcal{C} , and let $\tilde{\gamma} \geq 0$. Then the matrix

$$\mathbf{M} = \mathbf{I}_{|\mathcal{C}|} + \tilde{\gamma} \mathbf{V}_{\mathcal{C}} \mathbf{V}^{\top} \mathbf{V} \mathbf{V}_{\mathcal{C}}^{\top} \in \mathbb{R}^{|\mathcal{C}| \times |\mathcal{C}|}$$

*These authors contributed equally to this paper.

†Now at LinkedIn.

‡Corresponding author: zhipeng.wang@alumni.rice.edu

§Corresponding author: arazi@clemson.edu

is symmetric positive definite.

Since $\mathbf{V}^\top \mathbf{V} \in \mathbb{R}^{d \times d}$ is symmetric positive semidefinite, define

$$\mathbf{A} = \mathbf{V}_C \mathbf{V}^\top \mathbf{V} \mathbf{V}_C^\top \in \mathbb{R}^{|\mathcal{C}| \times |\mathcal{C}|}.$$

We first show that \mathbf{A} is positive semidefinite. For any $\mathbf{x} \in \mathbb{R}^{|\mathcal{C}|}$,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \mathbf{V}_C \mathbf{V}^\top \mathbf{V} \mathbf{V}_C^\top \mathbf{x} = (\mathbf{V}_C^\top \mathbf{x})^\top \mathbf{V}^\top \mathbf{V} (\mathbf{V}_C^\top \mathbf{x}).$$

Let $\mathbf{y} = \mathbf{V}_C^\top \mathbf{x} \in \mathbb{R}^d$. Then

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{y}^\top \mathbf{V}^\top \mathbf{V} \mathbf{y} = \|\mathbf{V} \mathbf{y}\|_2^2 \geq 0,$$

so \mathbf{A} is positive semidefinite.

Now consider

$$\mathbf{M} = \mathbf{I}_{|\mathcal{C}|} + \tilde{\gamma} \mathbf{A}.$$

The matrix \mathbf{M} is symmetric because both $\mathbf{I}_{|\mathcal{C}|}$ and \mathbf{A} are symmetric. For any nonzero $\mathbf{x} \in \mathbb{R}^{|\mathcal{C}|}$ we have

$$\mathbf{x}^\top \mathbf{M} \mathbf{x} = \mathbf{x}^\top \mathbf{I}_{|\mathcal{C}|} \mathbf{x} + \tilde{\gamma} \mathbf{x}^\top \mathbf{A} \mathbf{x} = \|\mathbf{x}\|_2^2 + \tilde{\gamma} \mathbf{x}^\top \mathbf{A} \mathbf{x}.$$

Since \mathbf{A} is positive semidefinite and $\tilde{\gamma} \geq 0$, we have $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$, so

$$\mathbf{x}^\top \mathbf{M} \mathbf{x} \geq \|\mathbf{x}\|_2^2 > 0$$

for all $\mathbf{x} \neq \mathbf{0}$. Therefore \mathbf{M} is positive definite.

Then we can apply Lemma 1, and we can conclude Γ is sumdoular. \square

B. Detail of Algorithm

In this appendix we describe how our pruning objective naturally leads to the OTPrune algorithm. The derivation begins directly with the token–token interaction matrix and proceeds to the incremental update rules used in Algorithm 1, without introducing additional notation beyond what is necessary.

Given the vision token matrix $\mathbf{V} \in \mathbb{R}^{m \times d}$, we first construct a token–token interaction matrix

$$\tilde{\mathbf{V}} = \mathbf{V} \mathbf{V}^\top \in \mathbb{R}^{m \times m}.$$

The i -th row of this matrix,

$$\mathbf{w}_i^\top = \tilde{\mathbf{V}}[i, :],$$

captures how token i interacts with all other tokens. For a subset $\mathcal{C} \subseteq [m]$, we use

$$\tilde{\mathbf{V}}_C = \tilde{\mathbf{V}}[\mathcal{C}, :]$$

to denote the corresponding row submatrix.

Our pruning objective is the log-determinant

$$\Gamma(\mathcal{C}) = \log \det \left(\mathbf{I} + \tilde{\gamma} \tilde{\mathbf{V}}_C \tilde{\mathbf{V}}_C^\top \right), \quad \tilde{\gamma} > 0.$$

This quantity becomes large when the selected interaction vectors $\{\mathbf{w}_i : i \in \mathcal{C}\}$ span a diverse subspace. It is convenient to define the kernel

$$\mathbf{K} = \mathbf{I} + \tilde{\gamma} \tilde{\mathbf{V}} \tilde{\mathbf{V}}^\top, \quad \text{so that} \quad K_{ij} = \delta_{ij} + \tilde{\gamma} \langle \mathbf{w}_i, \mathbf{w}_j \rangle.$$

Here, the Kronecker delta δ_{ij} equals 1 if $i = j$ and 0 otherwise. Then $\Gamma(\mathcal{C}) = \log \det(K_C)$, where K_C is the principal submatrix indexed by \mathcal{C} .

Greedy marginal gain. When considering a new token $j \notin \mathcal{C}$, the determinant of the enlarged matrix can be written using the classical one-step determinant identity:

$$\det(K_{\mathcal{C} \cup \{j\}}) = \det(K_C) \cdot \left(K_{jj} - K_{Cj}^\top K_C^{-1} K_{Cj} \right).$$

Thus the improvement in the objective is

$$\Gamma(\mathcal{C} \cup \{j\}) - \Gamma(\mathcal{C}) = \log d_j^2, \quad d_j^2 = K_{jj} - K_{Cj}^\top K_C^{-1} K_{Cj}.$$

The greedy algorithm simply selects the token with the largest d_j^2 .

Cholesky representation. Let the Cholesky factor of the current kernel submatrix be

$$K_C = LL^\top.$$

Define the coefficient vector

$$\mathbf{c}_j = L^{-1} K_{Cj}.$$

Because $K_C^{-1} = (L^\top)^{-1} L^{-1}$, the quadratic form simplifies to

$$K_{Cj}^\top K_C^{-1} K_{Cj} = \|\mathbf{c}_j\|_2^2,$$

so the greedy score becomes

$$d_j^2 = K_{jj} - \|\mathbf{c}_j\|_2^2.$$

Thus maintaining the greedy scores reduces to tracking \mathbf{c}_i and d_i^2 for all unselected tokens.

Adding a token. Suppose token j is selected. The Cholesky factor of the enlarged set $\mathcal{C} \cup \{j\}$ takes the form

$$L' = \begin{pmatrix} L & 0 \\ \mathbf{c}_j^\top & d_j \end{pmatrix}.$$

For any remaining token $i \notin \mathcal{C} \cup \{j\}$, let the updated coefficient vector be

$$\mathbf{c}'_i = \begin{bmatrix} \mathbf{c}_i \\ e_i \end{bmatrix}.$$

Plugging this into $L'c'_i = K_{\mathcal{C} \cup \{j\}, i}$ yields

$$e_i = \frac{K_{ji} - \langle c_j, c_i \rangle}{d_j}.$$

Once this new coordinate is known, the updated greedy score becomes

$$d_i^2 = K_{ii} - \|c'_i\|^2 = d_i^2 - e_i^2.$$

In general, the update takes the form

$$e_i = \frac{K_{ji} - \langle c_j, c_i \rangle}{d_j}.$$

For our kernel $K_{ji} = \delta_{ji} + \tilde{\gamma} \langle w_j, w_i \rangle$. Since we only update $i \neq j$ (token j has just been selected), the Kronecker delta term vanishes and we obtain

$$e_i = \frac{\tilde{\gamma} \langle w_j, w_i \rangle - \langle c_j, c_i \rangle}{d_j}.$$

This expression exactly the operations executed in Algorithm 1.

Initialization and complexity. With $\mathcal{C} = \emptyset$, all coefficient vectors are empty and the initial scores reduce to

$$d_i^2 = K_{ii} = 1 + \tilde{\gamma} \|w_i\|^2.$$

At iteration t , each of the $(m - t)$ remaining items requires an $O(t)$ inner product $\langle c_j, c_i \rangle$, so the overall cost of the iteration is

$$O((m - t)t).$$

Summing over $t = 1, \dots, k$ gives

$$\sum_{t=1}^k O((m - t)t) = O(mk^2),$$

which is the total complexity of selecting k tokens.

This completes the derivation of *OTPrune* directly from the log-determinant objective.

C. Discussion on Non-uniform Weighting Strategies

In our current framework, we adopt a *uniform weight assumption* for the source distribution P , where each token $v_i \in \mathbf{V}$ is assigned an equal mass of $1/M$. While this simplification allows for a highly efficient and training-free selection process, it treats all visual regions as equally significant, which may not fully capture the semantic variance inherent in complex images.

Potential for Adaptive Importance. Theoretically, the optimal transport formulation in *OTPrune* can be generalized

to accommodate *non-uniform weights* by assigning an importance score w_i to each token, such that $\sum w_i = 1$. As suggested during the rebuttal, such weights could be derived from internal model signals to better reflect semantic density:

- *Attention-based Priors:* Leveraging the cumulative attention weights from preceding transformer layers to identify patches that the model naturally prioritizes.
- *Task-specific Significance:* Utilizing cross-modal alignment scores to assign higher weights to visual tokens that are most relevant to the input text prompt.

Future Challenges. Transitioning to a non-uniform weighting strategy transforms *OTPrune* into a more complex joint optimization problem involving both subset selection and mass allocation. We consider the exploration of *soft-weighted OTPrune*, where token weights are dynamically adjusted during the pruning process—as a promising direction for further improving performance at extremely low pruning ratios.

References

- [1] Francis Bach. Learning with submodular functions: A convex optimization perspective. *arXiv preprint arXiv:1111.6453*, 2011. 1
- [2] Andreas Krause and Carlos Guestrin. Near-optimal observation selection using submodular functions. In *AAAI*, pages 1650–1654, 2007.
- [3] Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2-3):123–286, 2012. 1