Correlation-Aware Example Selection for In-Context Learning with Nonsymmetric Determinantal Point Processes

Qiunan Du¹, Zhiliang Tian^{1*}, Zhen Huang^{1*}, Kailun Bian¹, Tianlun Liu¹, Zhaoning Zhang¹, Xinwang Liu¹, Feng Liu², Dongsheng Li¹

¹College of Computer Science and Technology, National University of Defense Technology

²National University of Defense Technology

 $\{ \hbox{duqiunan, tianzhiliang, huangzhen, ltlun, zhangzhaoning, xinwangliu, richardlf} \} @nudt.edu.cn, \\ klbian@email.ncu.edu.cn$

Abstract

LLMs with in-context learning (ICL) obtain remarkable performance but are sensitive to the quality of ICL examples. Prior works on ICL example selection explored unsupervised heuristic methods and supervised LLM-based methods, but they typically focus on the selection of individual examples and ignore correlations among examples. Researchers use the determinantal point process (DPP) to model negative correlations among examples to select diverse examples. However, the DPP fails to model positive correlations among examples, while ICL still requires the positive correlations of examples to ensure the consistency of examples, which provides a clear instruction for LLMs. In this paper, we propose an ICL example selection method based on the nonsymmetric determinantal point process (NDPP) to capture positive and negative correlations, considering both the diversity and the relevance among ICL examples. Specifically, we optimize NDPP via kernel decompositionbased MLE to fit a constructed pseudo-labeled dataset, where we also propose a low-rank decomposition to reduce the computational cost. Further, we perform query-aware kernel adaptation on our NDPP to customize the input query, and we select examples via a MAP inference based on the adapted NDPP. Experimental results show our model outperforms strong baselines in ICL example selection.

1 Introduction

Large language models (LLMs) show good performance through in-context learning (ICL) (Brown et al., 2020; Wei et al., 2022b,a; Wen et al., 2024; Pan et al., 2024; Tian et al., 2025; Guo et al., 2024). ICL typically uses an example set and a task-specific instruction (with the user's query) as a prompt and feeds the prompt into LLMs. ICL allows LLMs to perform tasks by observing a series

of examples without the need to update parameters. However, the performance of ICL is sensitive to the selection of examples (Liu et al., 2022; Zhang et al., 2022; Min et al., 2022; An et al., 2023). Recent works (Lu et al., 2022; Cheng et al., 2023) also show that different example sets exhibit significant differences in performance, thus being crucial for exploiting the ICL capabilities of LLMs.

To select suitable examples for ICL, researchers propose various context-dependent heuristic methods, where they select examples according to the examples' entropy (Lu et al., 2022), complexity (Fu et al., 2022), perplexity (Gonen et al., 2023), and diversity (Li and Qiu, 2023). These methods outperform random selection, but these methods ignore characteristics of the specific input queries and thus cannot customize the ICL example set for the input queries. To consider the query, researchers propose context-aware methods to retrieve similar examples for ICL (Liu et al., 2022; Agrawal et al., 2023; Hongjin et al., 2022). They use off-the-shelf retrievers such as BM25 (Robertson et al., 2009) or SBERT (Reimers and Gurevych, 2019) to select examples based on their textual or semantic similarity to the query. When applying LLMs to specific tasks, they cannot customize the example selection of ICL for the given task since the ICL example selector (i.e., retriever) is not learnable and cannot learn to tailor to the task-specific data.

To leverage task supervision, some recent works (Rubin et al., 2022; Cheng et al., 2023; Li et al., 2023; Xiong et al., 2024) use LLMs' feedback as the task-specific supervisory signal to train the ICL example selectors (i.e., retriever), where the signal is used to rank and label examples. In these methods, the retrievers learn the LLMs' preference for examples in different tasks and adaptively select examples for each task. However, they typically focus on the selection of each individual example, ignoring the correlations (i.e., inter-relationships) among a set of ICL examples.

^{*} Corresponding Authors

To consider the correlations among examples for ICL, researchers (Levy et al., 2023; Ye et al., 2023a; Yang et al., 2023) propose to use the determinantal point process (DPP) (Kulesza and Taskar, 2012) to select examples by balancing the relevance to input queries and the diversity among examples. They model the relevance to input queries by similarity between queries and examples, and they model the diversity among examples since DPP's kernel matrix L models the negative correlation of data points. However, DPP's kernel matrix L is a symmetric positive semi-definite (PSD) matrix. L restricts DPP can only model negative correlation ¹ among examples rather than positive correlation. It results in DPP ignoring the relevance among candidate examples.

We argue that ICL example selection should not only consider the *relevance to input queries* and the *diversity among examples*, but also cater to the *relevance among examples*. Ensuring the consistency of ICL examples contributes to providing clear instructions to guide LLMs (Liu et al., 2024a). ²

In this paper, we propose an ICL example selection method for LLM based on the nonsymmetric determinantal point process model (NDPP), which considers the relevance to input queries, the diversity among ICL examples, and the relevance among ICL examples. NDPP's nonsymmetric property makes the selection consider relevance among ICL examples. Specifically, we construct an NDPP model with a kernel matrix to capture positive and negative correlations among ICL examples. In the training stage, we propose a kernel decompositionbased maximum likelihood estimation (KD-MLE) to train the NDPP by fitting the kernel matrix over our constructed pseudo-labeled datasets. To reduce the computational cost of KD-MLE, we propose a low-rank decomposition of the kernel matrix. In the inference stage, to consider the relevance to input queries, we propose a query-aware kernel adaptation, which adapts the trained NDPP to the given query by incorporating the embedding similarity between examples and queries into the kernel matrix. We finally perform maximal a posteriori (MAP) inference based on the adapted NDPP to

select the ICL example set for LLMs. Experiments show that our method exceeds baselines on five datasets, including open-domain QA, code generation, semantic parsing, and story generation tasks. Our code is released.³

Our contributions are: (1) We propose a novel ICL example selection framework based on NDPP, which captures positive and negative correlations among examples and models the composition of ICL examples to select suitable ICL examples for LLM. (2) We propose a query-aware kernel optimization to consider the similarity between queries and examples, which enables our method to select customized ICL example sets for different queries. (3) Experiments on five datasets show that our method achieves SOTA on ICL example selection.

2 Related Work

2.1 Example Selection for ICL

ICL example selection methods mainly have three categories: (1) In-context Insensitive Unsuper**vised Methods.** These approaches ignore the query information and task supervision. Researchers propose example selection methods based on complexity, entropy, diversity, and so on (Fu et al., 2022; Lu et al., 2022; Li and Qiu, 2023). (2) In-context **Sensitive Unsupervised Methods.** This category considers query information but ignores the task supervision. Researchers find that selecting different examples can reduce the redundancy of ICL example set (Liu et al., 2022; Agrawal et al., 2023; Hongjin et al., 2022). Wang et al. (2024a) propose a model-specific example selection method and Liu et al. (2024b) select examples with multiple levels of similarity to queries. (3) In-context Sensitive Supervised Methods. By introducing task supervision, these methods fine-tune ICL example selectors (i.e., retrievers). Many studies improved the quality of ICL examples by iteratively training retrievers (Rubin et al., 2022; Wang et al., 2024b; Li et al., 2023; Liu et al., 2024b). Xiong et al. (2024) further use chain-of-thought. (Levy et al., 2023; Yang et al., 2023; Ye et al., 2023b) use DPP to select diverse example sets. These works only consider relevance to input queries and diversity of examples, our model further considers relevance among examples.

¹In DPP, the correlation between examples i and j is expressed as $-\boldsymbol{L}_{ij}\boldsymbol{L}_{ji}$, where \boldsymbol{L} is the kernel matrix. Due to the symmetric property of PSD matrix, \boldsymbol{L}_{ij} and \boldsymbol{L}_{ji} are always equal, making the correlation $-\boldsymbol{L}_{ij}\boldsymbol{L}_{ji}$ always non-positive.

²The relevance and diversity are not conflicting since ICL needs multiple examples, where some of them may be diverse and others are relevant so as to provide a comprehensive and consistent instruction to LLMs.

³The implementation is available at: https://github.com/dqn1984/ICL-NDPP

2.2 DPP and Its Applications

(1) Theoretical studies on DPP. DPP has seen significant development. Johansson et al. (2023) proposed a semi-supervised k-DPP method. Grosse et al. (2024) used a greedy algorithm for k-DPP sampling. Okoth et al. (2022) propose LSMOEA-DPP and Ghilotti et al. (2024) propose Anisotropic DPP. (2) **Applications of DPP in AI.** DPP is widely used in AI applications, especially for tasks requiring diverse sets, e.g., neural network training (Sheikh et al., 2022), recommendation systems (Liu et al., 2024c), video analysis (Chen et al., 2023), and abstract summary (Shen et al., 2023). (3) Theoretical studies on DPP. Gartrell et al. (2019) propose NDPP, a nonsymmetric extension of DPP, which can model both positive and negative correlations among items. Gartrell et al. (2021) reduce NDPP's complexity. Han et al. (2022) propose a scalable sampling method for NDPP. Song et al. (2024) propose a fast dynamic algorithm for NDPP. While current works focus on the application of the DPP, we explore the application of the NDPP on ICL example selection. See more details of related work in App. A.

3 Preliminary

Nonsymmetric Determinantal Point Process.

NDPP is a probabilistic model to model correlations between items in a set (Gartrell et al., 2019). It models a finite ground set D with a kernel matrix L such that for any subset $E \subseteq D$, $P_L(E) \propto det(L_E)$, where L_E is the submatrix of L indexed by E. Given the kernel matrix L, the probability a subset E being selected from D is defined as:

$$P_{L}(E) = \frac{det(L_{E})}{det(L+I)}$$
 (1)

where I is a unit matrix. See App. B for more details of NDPP and ICL.

4 Method

4.1 Overview

To provide high-quality ICL examples for LLMs, we construct an ICL example selection framework based on the NDPP model, where the NDPP consists of a kernel matrix L to model correlations among examples. We construct a pseudo-labeled training set based on LLMs' feedback (§ 4.2), and use the pseudo-labeled training set to train the

NDPP model by kernel decomposition-based maximum likelihood estimation (MLE) (§ 4.3). In the inference stage, we perform query-aware kernel-adaptation on the trained NDPP model to consider the relevance to input queries, and select ICL examples based on the adapted model through a maximum a posteriori (MAP) inference (§ 4.4).

4.2 Example Subsets Pseudo-labeling via LLMs' Feedback

Since there is no ground truth of ICL example sets for each training instance, to train the NDPP model in § 4.3 by MLE, we collect the feedback signals from LLMs for scoring the example subsets to construct a pseudo training set.

Given a task, we construct a pseudo-labeled training set with three steps: (1) Candidate ex**ample retrieval.** For each instance (x_i, y_i) from our training set, we retrieve a candidate example set from the example pool D using the KNN retriever, which considers the embedding similarity between the instance and examples. From the retrieved candidate example set, we randomly sample N non-overlapping subsets, denoted as $\{E_{ij}\}_{i=1}^{N}$. (2) **Example subset scoring.** We measure the quality of each candidate example subset E_{ij} with a quality score s_{ij} , and the scores act as soft pseudo labels of the subsets. To obtain the quality score s_{ij} , we concatenate the query x_i and examples in the subset E_{ij} , and input the concatenation into an LLM to obtain the probability $P_{LLM}(y_i|E_{ii},x_i)$ of predicting the corresponding ground truth y_i of the test query x_i , which is formalized as: $s_{ij} = P_{LLM}(y_i|E_{ij},x_i)$. (3) **Pseudo** training set construction. We rank candidate example subsets based on the score s_{ij} , and select the top 10% high-scoring subsets for all instances to construct a pseudo-labeled training set $D_{train} = (E_i)_{i=1}^n$, where n is the subset number. D_{train} is used to train the NDPP model in (§ 4.3).

4.3 NDPP Model Optimization with Pseudo-labeled Example Subsets

To select high-quality ICL example sets, we train the NDPP model by kernel decomposition-based MLE, which allows the NDPP model to learn the kernel matrix of high-scoring example subsets from the pseudo-labeled training set. The process consists of three steps: (1) we first define the NDPP optimization objective, then (2) get the kernel decomposition for NDPP, and finally, (3) we optimize NDPP via the kernel decomposition-based MLE.

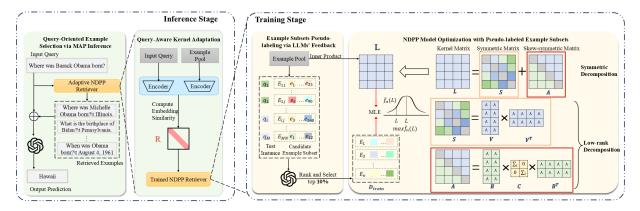


Figure 1: The overview of our framework. In the training stage, we construct a pseudo-labeled training set D_{train} based on LLMs' feedback (§ 4.2), and use D_{train} to optimize the kernel matrix \boldsymbol{L} of the NDPP model by kernel decomposition-based MLE (§ 4.3). In the inference stage, we perform query-aware kernel-adaptation on the trained NDPP model, and select ICL examples based on the adapted model through MAP inference (§ 4.4).

4.3.1 NDPP Optimization Objective: MLE with Kernel Matrix

To capture correlations among ICL examples, we optimize the kernel matrix of the ICL example set to fit the pseudo-labeled training set. The fitted kernel matrix represents the feature of high-scoring ICL example sets so that NDPP can select suitable examples with the fitted kernel matrix.

In the NDPP, recall that the probability of selecting a candidate example subset E_i from the example pool D is $P_{\boldsymbol{L}}(E_i) = \frac{det(\boldsymbol{L}_{E_i})}{det(\boldsymbol{L}+I)}$ (as shown in Eq. 1), where \boldsymbol{L} is the kernel matrix of D and \boldsymbol{L}_{E_i} is the submatrix of \boldsymbol{L} indexed by E_i . \boldsymbol{L} is constructed by computing the pairwise embedding similarity between two examples $\langle \boldsymbol{e}_i, \boldsymbol{e}_j \rangle$ in the example pool D, where $\boldsymbol{L}_{ij} = sim(\boldsymbol{e}_i, \boldsymbol{e}_j)$. Elements of \boldsymbol{L} show correlations among examples in D. Given different kernel matrices, NDPP selects different ICL example sets with the probability $P_{\boldsymbol{L}}(\cdot)$.

To select high-quality ICL example sets with NDPP, we aim to find a kernel matrix \boldsymbol{L} that maximizes the probability of selecting high-scoring ICL example subsets. To achieve it, we optimize the kernel matrix \boldsymbol{L} of the ICL example set to fit the pseudo-labeled training set $D_{train} = (E_i)_{i=1}^n$. Specifically, we optimize \boldsymbol{L} towards the log-likelihood on the training set D_{train} as,

$$\hat{f}_n(L) = \frac{1}{n} \sum_{i=1}^n log P_L(E_i)$$
 (2)

Because $P_{\boldsymbol{L}}(E_i) = \frac{det(\boldsymbol{L}_{E_i})}{det(\boldsymbol{L}+\boldsymbol{I})}$, we have:

$$\hat{f}_n(\mathbf{L}) = \frac{1}{n} \sum_{i=1}^n \log \det(\mathbf{L}_{E_i}) - \log \det(\mathbf{L} + \mathbf{I}) \quad (3)$$

The optimized kernel matrix \hat{L} is the kernel matrix that maximizes the Eq. 3, denoted as:

$$\hat{\boldsymbol{L}} = \arg\max_{\boldsymbol{L}} \hat{f}_n(\boldsymbol{L}) \tag{4}$$

The convexity analysis of Eq. 4 is provided in App. C. The optimized kernel matrix \hat{L} is the learnable optimal approximation of high-scoring ICL example subsets' kernel matrix, with its elements representing correlations among examples.

4.3.2 Kernel Decomposition of NDPP

To optimize the kernel matrix \boldsymbol{L} conveniently, we perform a two-step decomposition on the NDPP kernel matrix: we first perform symmetric decomposition on the kernel matrix, which enables NDPP to learn the positive and negative correlations among examples independently, and then perform a low-rank decomposition to reduce the computational cost. Details are as follows:

Symmetric decomposition. To distinguish the positive and negative correlations among examples (using NDPP's nonsymmetric property), we decompose the kernel matrix \boldsymbol{L} into the sum of a symmetric matrix \boldsymbol{S} and a skew-symmetric matrix \boldsymbol{A} as in Eq. 5, where \boldsymbol{A} and \boldsymbol{S} denote the positive and negative correlations, respectively.

Low-rank decomposition. To reduce the computational cost, inspired by Gartrell et al. (2021), we further perform a low-rank decomposition on the symmetric matrix S and the skew-symmetric matrix A as in Eq. 5, which converts the high-dimensional representation of the correlations into a low-dimensional representation.

$$L = S + A, S = VV^{T}, A = BCB^{T}$$
 (5)

 $oldsymbol{V}, oldsymbol{B} \in \mathbb{R}^{M imes K}$ are low-rank matrices of $oldsymbol{S}$ and $oldsymbol{A}$ respectively, where M is the example number in the example pool D and K is the rank of the kernel matrix $oldsymbol{L}$. $oldsymbol{V}$ and $oldsymbol{B}$ indicate the low-dimensional representation of the negative and positive correlations among examples, respectively. $oldsymbol{C} \in \mathbb{R}^{K imes K}$ is a block-diagonal matrix with diagonal blocks Σ_i of the form $egin{bmatrix} 0 & \lambda_i \\ -\lambda_i & 0 \end{bmatrix}$, where $\lambda_i > 0$. $oldsymbol{C}$ maintains the skew-symmetric property of $oldsymbol{A}$.

4.3.3 Kernel Decomposition-based MLE

We perform MLE to fit the kernel matrix L with its kernel decomposition form $L = VV^T + BCB^T$ obtained in the above step, where we also apply a regularization term to the log-likelihood.

Step 1: Kernel-decomposed MLE. When we optimize the kernel matrix L towards the MLE objective, we need to perform the decomposition of L to ensure that L captures both positive and negative correlations. We recall the log-likelihood of L (Eq. 3). Specifically, we use the decomposition form $L = VV^T + BCB^T$ in Eq. 5 to decompose L and L_{E_i} in the objective function (Eq. 3) to obtain the kernel-decomposed log-likelihood (Eq. 6),

$$\phi(\boldsymbol{V}, \boldsymbol{B}, \boldsymbol{C})$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log \det \left(\boldsymbol{V}_{E_i} \boldsymbol{V}_{E_i}^T + \boldsymbol{B}_{E_i} \boldsymbol{C} \boldsymbol{B}_{E_i}^T \right)$$

$$- \log \det \left(\boldsymbol{V} \boldsymbol{V}^T + \boldsymbol{B} \boldsymbol{C} \boldsymbol{B}^T + \boldsymbol{I} \right)$$
(6)

Eq. 6 allows us to optimize the log-likelihood with the decomposed components V, B, C. The matrices B and V can capture positive and negative correlations among examples, respectively. Note the second term det(L+I) in Eq. 3 requires calculation with complexity $O(M^3)$, while the kernel decomposition reduce the computational complexity of the second term in Eq. 6 to $O(MK^2 + nK^3)$. The running time linearly scales with the dataset size M, and we show the time cost in table 7.

Step 2: Regularized log-likelihood. To prevent overfitting, we define a regularization term as shown in Eq. 7. We perform L2 regularization for each row vector v_i and b_i of the matrices V and B separately, and use hyperparameters α and β to control the regularization strength of the matrices V and B, respectively. Besides, we define a weight parameter $\frac{1}{\gamma_i}$ to control the regularization strength for each row vector, where γ_i denotes the occurrences of the i_{th} element appears in D_{train} .

The regularization term is formally denoted as:

$$R(\boldsymbol{V}, \boldsymbol{B}) = -\alpha \sum_{i=1}^{M} \frac{1}{\gamma_i} \| \boldsymbol{v}_i \|_2^2 - \beta \sum_{i=1}^{M} \frac{1}{\gamma_i} \| \boldsymbol{b}_i \|_2^2 \quad (7)$$

Adding the regularization term (Eq. 7) to the kernel-decomposed log-likelihood (Eq. 6), we obtain the regularized log-likelihood (Eq. 8):

$$\phi(\mathbf{V}, \mathbf{B}, \mathbf{C})$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log \det \left(\mathbf{V}_{E_i} \mathbf{V}_{E_i}^T + \mathbf{B}_{E_i} \mathbf{C} \mathbf{B}_{E_i}^T \right)$$

$$- \log \det \left(\mathbf{V} \mathbf{V}^T + \mathbf{B} \mathbf{C} \mathbf{B}^T + I \right)$$

$$+ R(\mathbf{V}, \mathbf{B})$$
(8)

The practical optimization process of Eq. 8 is provided in App. D.

In summary of the processing of § 4.3, we first train the NDPP model on the pseudo-labeled training set D_{train} collected in § 4.2, where we optimize Eq. 8 to find the optimized kernel matrix (§ 4.3.1) \hat{L} through its kernel decomposition form (§ 4.3.2 and§ 4.3.3) as Eq. 5. Then, the optimized kernel matrix can assist the NDPP model to select high-quality ICL example sets.

4.4 ICL Example Selection via NDPP for LLMs Inference

In the inference stage, to provide customized high-quality ICL examples for different queries, we propose query-aware kernel adaptation to adapt the trained NDPP to specific input queries so as to select ICL examples. To achieve it, we adapt the NDPP to input queries by modeling the similarity between examples and queries (§ 4.4.1), and then select ICL examples by maximum a posteriori (MAP) inference using the adapted NDPP (§ 4.4.2). The above operations consider both the relevance to input queries and the relevance among examples.

4.4.1 Adapting NDPP to Input Queries

To adapt NDPP to input queries, we update its kernel matrix by introducing the similarity between examples and input queries into the kernel matrix.

For each query, we update the kernel matrix with three steps: (1) **Similarity Score Computation.** We encode the query x via a query encoder $E_Q(\cdot)$ and encode the example e_i via an example encoder $E_P(\cdot)$. We obtain the similarity score r_i via the inner product of their encoder outputs: $r_i = \sin(x, e_i) = E_Q(x)^T E_P(e_i)$. (2) **Similarity Matrix Construction.** Using similarity scores

 $r = [r_1, r_2, ..., r_M]$ for all M examples in the example pool D, we construct a diagonal similarity matrix $R \in \mathbb{R}^{M \times M}$: R = Diag(r), where $Diag(\cdot)$ is the diagonal matrix operator. The diagonal of R consists of r, while all off-diagonal elements are 0. (3) **Kernel Matrix Adaptation.** We adapt the optimized kernel matrix to the given input query by incorporating the above similarity matrix R with the optimized kernel matrix \hat{L} obtained in § 4.3. That is, we obtain the adapted kernel matrix L' as: $L' = R \cdot \hat{L} \cdot R$.

4.4.2 Query-Oriented Example Selection via MAP Inference

To select the ICL example set for queries with the adapted NDPP, rather than selecting the most relevant k examples (Rubin et al., 2022; Wang et al., 2024b), we conduct the MAP inference, the standard subset sampling method for NDPP when the application requires a single output set(Gartrell et al., 2021), to select examples one by one from the example pool D via a greedy algorithm. The goal of MAP inference is to select the high-quality ICL example set S_{map} of size k from D for the current query. In the adapted NDPP, given the kernel matrix L', S_{map} is the example subset of size k from D that maximizes $P_{L'}(S)$ among all possible subsets S of size k. Recall that the probability $P_{L'}(S)$ is proportional to the determinant of the sub-kernel matrix L'_S , S_{map} is the example subset of D that maximizes $det(\mathbf{L}'_S)$ among all subsets S of size k. Formally, we define the MAP inference of the example selection with adapted NDPP as:

$$S_{map} = \underset{S \subseteq D, |S|=k}{\arg\max} \log \det(\mathbf{L}_S')$$
 (9)

However, the MAP inference above has been proved to be NP-hard⁴ (Ko et al., 1995; Kulesza and Taskar, 2012). To reduce the computational cost, a common approach is to approximate the MAP inference using greedy algorithms (Nemhauser et al., 1978; Gillenwater et al., 2012; Chen et al., 2018). To reduce the cost, we first select a candidate example set Z, |Z| = m, m < M with a KNN retriever to reduce the size of candidate examples. Then,

following Gartrell et al. (2021), we approximate MAP inference using a greedy algorithm with the complexity of $O(mKk+mK^2)$: starting from an empty set S_{map} , we iteratively select examples one by one until we obtained k examples, approximating the global optimum by solving local optima at each iteration. At each iteration, for all examples i in the candidate example set Z that are not included in S_{map} , we compute the increment of the log-determinant $logdet(\cdot)$ of the sub-kernel matrix $L'_{S_{map}}$ after adding example i to the set S_{map} . We select the example j with the largest increment as the local optimum and add it into S_{map} :

$$j = \underset{i \in Z \backslash S_{map}}{\operatorname{arg \, max}} \log \det \left(\mathbf{L}'_{S_{map} \bigcup \{i\}} \right)$$

$$-\log \det \left(\mathbf{L}'_{S_{map}} \right)$$
(10)

App. E provides a proof of a lower bound on the approximation quality of the greedy algorithm. Finally, we concatenate the query and the ICL example set S_{map} as the input prompt of LLMs.

5 Experiments

5.1 Experimental Settings

Dataset. Following (Ye et al., 2023b; Li et al., 2023), we use five datasets: (1) GeoQuery (Shaw et al., 2020) has 880 geography questions. (2) NL2Bash (Lin et al., 2018) contains 9k Bash command pairs. (3) MTOP (Li et al., 2020) is a multilingual parsing dataset with 6 languages. (4) WebQs (Berant et al., 2013) covers 6,642 QA pairs using Freebase. (5) Roc End (Mostafazadeh et al., 2016) is a corpus with 100k stories.

Metrics. Following (Ye et al., 2023b; Li et al., 2023), we use those metrics: (1) Exact Match (EM) (Rajpurkar et al., 2016) for GeoQuery, MTOP, and WebQs to assess the accuracy of the generated output. (2) BLEU-1 (Papineni et al., 2002) for Roc Ending to evaluate alignment in story generation. (3) BLEU-4 (Papineni et al., 2002) for NL2Bash to capture longer sequence structure in command generation.

Baselines. We compare with two types of methods: (1) **Unsupervised Methods:** Random, which randomly selects non-repeating ICL examples from the example pool. BM25 (Robertson et al., 2009), which extends TF-IDF to rank relevant examples for the test input and select the top-k highest scoring ICL examples for each test input. (2) **Supervised Methods:** EPR (Rubin et al., 2022), which uses the LLM itself as a scoring model to retrieve

 $^{^4}$ Such MAP inference requires finding all subsets S, |S| = k of the example pool D, |D| = M and computing their determinants. The example pool D, |D| = M has C(M,k) subsets S, |S| = k in total, and the computational complexity of each subset determinant is $O(k^3)$. The cost of the MAP inference is $O(M^k \cdot k^3)$ in total, which is unaffordable as the size of the example pool D increases. And the function $logdet(\boldsymbol{L}_S')$ is proved to be submodular, and the unconstrained optimization problem for submodular is NP-hard.

| Model | Method | GeoQuery (EM) | MTOP (EM) | NL2Bash (BLEU-4) | WebQs (EM) | RocEnd (BLEU-1) |
|----------------|--------|------------------|--------------|---------------------|---------------|--------------------|
| GPT-Neo (2.7B) | Random | 33.57 | 0.67 | 34.35 | 4.87 | 57.58 |
| | BM25 | 62.86 | 53.24 | 58.98 | 16.68 | 58.65 |
| | EPR | 71.07 | 60.36 | 56.82 | 17.91 | 59.12 |
| | CEIL* | 70.71 | 63.40 | 53.66 | 17.08 | 59.72 |
| | TTF* | 68.93 | 54.05 | 56.11 | 16.14 | / |
| | Our | 73.21 | 65.37 | 61.01 | 18.90 | 60.33 |
| GPT-4 | Random | 71.43 | 21.48 | 67.45 | 34.49 | 58.34 |
| | EPR | 88.93 | 78.61 | 73.63 | 50.32 | 54.70 |
| | CEIL | 91.07 | 78.70 | 73.95 | 46.75 | 56.24 |
| | Our | 91.43 | 79.02 | 73.96 | 52.95 | 62.81 |

Table 1: ICL example selection experiment results. "/" indicates that the method is not open source and does not give results of the dataset in the corresponding paper and "Bold" indicates optimal results. All results are averaged over 3 runs. We reference results from the previous work (Liu et al., 2024b), marked by *. Our improvements are significant under the t-test with p < 0.05 (See details in App. H).

| Settings | GeoQuery | MTOP | NL2Bash | WebQs | RocEnd |
|---|--------------------------------|--------------------------------|--------------------------------|----------------------------------|--------------------------------|
| | (EM) | (EM) | (BLEU-4) | (EM) | (BLEU-1) |
| Ours(Full Model) w/o Scoring w/o Regularization w/o Adaptation | 73.21 72.36 71.43 71.64 | 65.37 65.19 65.28 65.28 | 61.01 59.32 60.25 59.56 | 18.90 17.91 18.75 18.45 | 60.33 59.09 59.94 60.33 |

Table 2: Ablation study. w/o Scoring: remove the LLM scoring when construct the training set; w/o Regularization: remove the regularization term in the log-likelihood; w/o Adaptation: remove query-aware kernel adaptation on the trained NDPP.

good ICL examples. CEIL (Ye et al., 2023b) models ICL example sets with DPP and trains DPP by contrastive learning. TTF (Liu et al., 2024b) finetunes the ICL example selector with labeled data, adding task-specific modules.

See details of datasets, metrics, baselines, and implementation in App. F.

5.2 Overall Performance

Table 1 shows the overall results of ICL example selection methods across five datasets. Notably, while prior studies (Ye et al., 2023a) primarily focus on smaller models like GPT-Neo (2.7B), we extend the evaluation to the SOTA LLM GPT-4⁵. The results demonstrate that our method outperforms all baseline methods on both GPT-neo (2.7B) and GPT-4 models , indicating the effectiveness of NDPP for ICL example selection.

Compared to random selection, our method shows over 20% average improvement on both models. All designed selection methods outperform random selection (except GPT-4 on RocEnd), highlighting the value of careful example selection.

We observe that the improvement of our method is more pronounced on GPT-neo (2.7B) compared to GPT-4, likely due to the latter's inherently stronger inference capability. This finding is consistent with previous research (Zhang et al., 2022). Notably, on Geoquery, Mtop, and RocEnd, our method on GPT-neo (2.7B) outperforms random example selection on GPT-4, demonstrating the effectiveness of our method in enhancing the ICL capability of LLMs. Furthermore, our method consistently outperforms CEIL on all datasets, suggesting the benefits of capturing positive correlations among examples for ICL example selection.

5.3 Ablation Studies

Table 2 presents results of the ablation study. Our complete model performs excellently across all five datasets, and removing any single module leads to a decrease in performance, validating the effectiveness of each component. Specifically: (1) w/o Scoring: We remove the step of scoring with LLM and instead use all the example subsets as the training set. We observe that although performance slightly declined, our model still maintains relatively good performance on some tasks. This suggests that our model is still able to model correlations among examples to some extent, but is disturbed by noise in low-scoring ICL example subsets. (2) w/o Regu-

⁵Due to the limitations of black-box models like GPT-4 (which only expose log probabilities for the first five tokens), our framework cannot directly construct pseudo-labeled training sets based on full token probabilities. To address this, we transfer the retriever trained on GPT-Neo (2.7B) directly to GPT-4 for ICL example selection.

larization: We removed the regularization term in Eq. 8, and the performance of our model deteriorates on certain tasks. Without regularization, our model exhibits a tendency to overfit, which results in a decrease in generalization ability on test data. (3) w/o Adaptation: We remove the query-aware kernel adaptation and observe a performance drop, which demonstrates the importance of considering the relevance between queries and examples.

5.4 Analysis studies

Performance Over Different Example Order.

Previous work (Lu et al., 2022) showed that ICL is sensitive to the order of examples when selecting examples randomly. We conduct experiments to investigate the effect of ordering on ICL examples retrieved by our method. Specifically, we provide 8 examples with 10 different random orderings for each dataset. We present the best (Best Random-Order) and worst (Worst Random-Order) results and the variance of the results over 10 runs. Results in Table 3 show that performance fluctuates somewhat across different orderings, but the variation is relatively small and within a controllable range. This suggests that although examples' order does have some impact on the performance of our model, the effect is limited. This finding is consistent with previous research (Li and Qiu, 2023), which indicates that high-quality examples can reduce ICL sensitivity to the order of examples.

| | GeoQuery | MTOP | WebQs | RocEnd |
|--------------------|----------|-------|-------|--------|
| Best Random-Order | 69.29 | 62.64 | 14.86 | 59.50 |
| Worst Random-Order | 66.43 | 61.48 | 13.24 | 58.10 |
| VAR | 0.78 | 0.13 | 0.21 | 0.19 |

Table 3: Performance over different example orders.

Performance Over Different Example Numbers. Many LLMs are constrained by limited input lengths, which restricts the maximum number of ICL examples that can be provided. To analyze the impact of example quantity on ICL performance, we compared three methods across four tasks, and the results are shown in Figure 2. Our key observations are as follows: Increasing the number of examples enhances ICL performance, as additional examples enable LLMs to better understand the task objectives and output patterns.

Performance Over Different Example Selection Methods. We compare two ICL example selec-

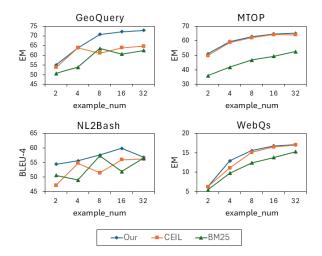


Figure 2: Performance over different example numbers.

| Method | GeoQuery | MTOP | NL2Bash | WebQs | RocEnd |
|--------|----------|-------|---------|-------|--------|
| MCMC | 53.21 | 62.19 | 52.83 | 15.11 | 57.09 |
| Our | 73.21 | 65.37 | 61.01 | 18.90 | 60.33 |

Table 4: Performance over different example selection methods in the inference stage.

tion methods (i.e., MCMC sample and greedy (our)) in the inference stage on the GPT-Neo(2.7B) model. MCMC sample (Gartrell et al., 2021) explores the subset space through random walking and probabilistic acceptance mechanisms. Results in Table 4 show that the quality of the ICL example set selected by our greedy algorithm is better than the set selected by the MCMC sample method, indicating the effectiveness of our greedy algorithm.

Performance Over Different Regularization Hyperparameters. We analyze varying regularization parameters α and β across datasets. The experimental results in Table 6 show that: (1) the combination of α , β we used performs best. (2) The RocEnd dataset is almost unaffected by the regularization parameter. This may be due to the fact that the RocEnd dataset has a much larger amount of data than the other datasets, and the risk of model overfitting is very low.

| Method | GeoQuery | MTOP | NL2Bash | WebQs | RocEnd |
|--------|--------------|--------------|--------------|--------------|--------------|
| CEIL | 82.14 | 81.12 | 72.93 | 30.59 | 59.50 |
| Our | 87.86 | 81.79 | 73.33 | 31.59 | 60.03 |

Table 5: Experimental results on Gemini-2.0-flash.

Generalization on other LLMs. To validate the effectiveness of our method on LLMs beyond the GPT family, we conduct experiments on the

Gemini-2.0-flash model. We compared our method with the best-performing baseline CEIL, and the results show that our method outperforms the baseline on all datasets, demonstrating the consistent superiority of our method across different models. The results are shown in Table 5.

Case Study. For an intuitive grasp of the effectiveness of our method, we present some cases on the RocEnd task. In Table 10, we compare prompts from Random, BM25, EPR, CEIL, and our method. Random-selected examples are irrelevant to the query; BM25-selected examples focus on keyword matching, but include redundant content such as the Niagara Falls honeymoon and the Hawaii honeymoon. EPR-selected examples do not ensure emotional consistency and deviate from the topic. CEIL-selected examples are more diverse but mixed with negative events, resulting in completely opposite outcomes. Compared to baselines, our method considers both the relevance between examples and queries, as well as the relevance and diversity among examples, ensuring that the provided examples maintain relevance to the query while comprehensively covering subsequent story generation patterns and maintaining thematic and emotional consistency, thereby providing appropriate guidance for the LLM to generate the target content. See App. G for more detailed analysis.

6 Conclusion

In summary, we proposed an NDPP-based framework for ICL example selection. Our framework first constructs a pseudo-labeled training set based on LLM feedback, and then uses the set to train the NDPP model by kernel decomposition-based MLE. Finally, in the inference stage, we perform query adaptation on the NDPP model, followed by MAP inference to select suitable and customized ICL example sets for different queries. Our experiments on five datasets across four domains show that our framework achieves SOTA performance in ICL example selection.

Limitations

The pseudo-labeled training dataset we construct relies on LLM feedback, which may be subject to inherent biases within the LLM. To address this limitation, future work could explore integrating fairness-aware mechanisms into the LLM feedback process, such as debiasing techniques, fairness constraints, or adversarial training, to mitigate potential biases.

Our framework constructs pseudo-labeled datasets based on token probabilities from LLM feedback, which inherently limits its compatibility with black-box models (e.g., GPT-4), as they only expose log probabilities for the top five tokens. However, our experiments demonstrate that a retriever trained on white-box models (e.g., GPT-Neo) can be effectively transferred to black-box models, achieving competitive performance. In future work, we plan to explore alternative approaches for constructing pseudo-labeled datasets that are universally applicable, including black-box LLMs.

Ethical Considerations

Privacy: This study utilizes only publicly available open datasets for all experiments, ensuring that no private or sensitive data is involved. The method does not collect, store, or process any personal information, thereby posing no risk to individual privacy.

Human Resources: Our research does not involve any manual annotation or human subject participation. All training and evaluation processes are automated, eliminating concerns regarding labor exploitation, unfair compensation, or excessive workloads.

Methodological Application: We believe that this study contributes intellectual value to the dependable application of retrieval-based in-context learning in the field of NLP, with potential broader implications for tasks in other areas. Potential biases present in the LLM may propagate through the retrieval process. We encourage users to critically evaluate the outputs and consider incorporating debiasing techniques when deploying such systems in sensitive applications.

Acknowledgements

This work is supported by the following foundations: the National Natural Science Foundation of China (NSFC) under Grant No.62376284, No.62306330, No.62202487, No.U22B2005, the National Science Foundation for Distinguished Young Scholars under Grant No.62025208, No.62325604, No.62125604, and the Young Elite Scientist Sponsorship Program by CAST under Grant No.YESS20230367.

References

- Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. Incontext examples selection for machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873.
- Shengnan An, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Jian-Guang Lou, and Dongmei Zhang. 2023. How do in-context examples affect compositional generalization? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11027–11052.
- Maurice Stevenson Bartlett. 1937. Properties of sufficiency and statistical tests. *Proceedings of the Royal Society of London. Series A-Mathematical and Physical Sciences*, 160(901):268–282.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1533–1544.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Laming Chen, Guoxin Zhang, and Eric Zhou. 2018. Fast greedy map inference for determinantal point process to improve recommendation diversity. *Advances in Neural Information Processing Systems*, 31
- Xiwen Chen, Huayu Li, Rahul Amin, and Abolfazl Razi. 2023. Rd-dpp: Rate-distortion theory meets determinantal point process to diversify learning data samples. *arXiv* preprint arXiv:2304.04137.
- Daixuan Cheng, Shaohan Huang, Junyu Bi, Yuefeng Zhan, Jianfeng Liu, Yujing Wang, Hao Sun, Furu Wei, Weiwei Deng, and Qi Zhang. 2023. Uprise: Universal prompt retrieval for improving zero-shot evaluation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12318–12337.
- Wei Duan, Junyu Xuan, Maoying Qiao, and Jie Lu. 2022. Learning from the dark: boosting graph convolutional neural networks with diverse negative samples. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6550–6558.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2022. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*.

- Mike Gartrell, Victor-Emmanuel Brunel, Elvis Dohmatob, and Syrine Krichene. 2019. Learning nonsymmetric determinantal point processes. *Advances in Neural Information Processing Systems*, 32.
- Mike Gartrell, Insu Han, Elvis Dohmatob, Jennifer Gillenwater, and Victor-Emmanuel Brunel. 2021. Scalable learning and {map} inference for nonsymmetric determinantal point processes. In *International Conference on Learning Representations*.
- Lorenzo Ghilotti, Mario Beraha, and Alessandra Guglielmi. 2024. Bayesian clustering of high-dimensional data via latent repulsive mixtures. *Biometrika*, page asae059.
- Jennifer Gillenwater, Alex Kulesza, and Ben Taskar. 2012. Near-optimal map inference for determinantal point processes. *Advances in Neural Information Processing Systems*, 25.
- Hila Gonen, Srini Iyer, Terra Blevins, Noah A Smith, and Luke Zettlemoyer. 2023. Demystifying prompts in language models via perplexity estimation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10136–10148.
- Julia Grosse, Rahel Fischer, Roman Garnett, and Philipp Hennig. 2024. A greedy approximation for kdeterminantal point processes. In *International Con*ference on Artificial Intelligence and Statistics, pages 3052–3060. PMLR.
- Yuxuan Guo, Zhiliang Tian, Yiping Song, Tianlun Liu, Liang Ding, and Dongsheng Li. 2024. Context-aware watermark with semantic balanced green-red lists for large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 22633–22646.
- Insu Han, Mike Gartrell, Jennifer Gillenwater, Elvis Dohmatob, and Amin Karbasi. 2022. Scalable sampling for nonsymmetric determinantal point processes. *arXiv preprint arXiv:2201.08417*.
- SU Hongjin, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. 2022. Selective annotation makes language models better fewshot learners. In *The Eleventh International Conference on Learning Representations*.
- Simon Johansson, Ola Engkvist, Morteza Haghir Chehreghani, and Alexander Schliep. 2023. Diverse data expansion with semi-supervised k-determinantal point processes. In 2023 IEEE International Conference on Big Data (BigData), pages 5260–5265. IEEE.
- Chun-Wa Ko, Jon Lee, and Maurice Queyranne. 1995. An exact algorithm for maximum entropy sampling. *Operations Research*, 43(4):684–691.
- Alex Kulesza and Ben Taskar. 2012. Determinantal point processes for machine learning.

- Itay Levy, Ben Bogin, and Jonathan Berant. 2023. Diverse demonstrations improve in-context compositional generalization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1401–1422.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2020. Mtop: A comprehensive multilingual task-oriented semantic parsing benchmark. *arXiv preprint arXiv:2008.09335*.
- Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4644–4668.
- Xiaonan Li and Xipeng Qiu. 2023. Finding support examples for in-context learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6219–6235, Singapore. Association for Computational Linguistics.
- Xi Victoria Lin, Chenglong Wang, Luke Zettlemoyer, and Michael D Ernst. 2018. Nl2bash: A corpus and semantic parser for natural language interface to the linux operating system. *arXiv* preprint *arXiv*:1802.08979.
- Haoyu Liu, Jianfeng Liu, Shaohan Huang, Yuefeng Zhan, Hao Sun, Weiwei Deng, Furu Wei, and Qi Zhang. 2024a. se2: Sequential example selection for in-context learning. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 5262–5284.
- Hui Liu, Wenya Wang, Hao Sun, Chris Xing Tian, Chenqi Kong, Xin Dong, and Haoliang Li. 2024b. Unraveling the mechanics of learning-based demonstration selection for in-context learning. *arXiv* preprint arXiv:2406.11890.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for gpt-3? In Proceedings of Deep Learning Inside Out (Dee-LIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, pages 100–114.
- Yuli Liu, Christian Walder, and Lexing Xie. 2024c. Learning k-determinantal point processes for personalized ranking. In 2024 IEEE 40th International Conference on Data Engineering (ICDE), pages 1036–1049. IEEE.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming fewshot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098.

- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022. Metaicl: Learning to learn in context. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2791–2809.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849.
- George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. 1978. An analysis of approximations for maximizing submodular set functions—i. *Mathematical programming*, 14:265–294.
- Michael Aggrey Okoth, Ronghua Shang, Licheng Jiao, Jehangir Arshad, Ateeq Ur Rehman, and Habib Hamam. 2022. A large scale evolutionary algorithm based on determinantal point processes for large scale multi-objective optimization problems. *Electronics*, 11(20):3317.
- Shilong Pan, Zhiliang Tian, Liang Ding, Haoqi Zheng, Zhen Huang, Zhihua Wen, and Dongsheng Li. 2024. POMP: Probability-driven meta-graph prompter for LLMs in low-resource unsupervised neural machine translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9976–9992, Bangkok, Thailand. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671.
- Peter Shaw, Ming-Wei Chang, Panupong Pasupat, and Kristina Toutanova. 2020. Compositional generalization and natural language variation: Can a semantic parsing approach handle both? *arXiv preprint arXiv:2010.12725*.
- Hassam Sheikh, Kizza Frisbee, and Mariano Phielipp. 2022. Dns: Determinantal point process based neural network sampler for ensemble reinforcement learning. In *International Conference on Machine Learning*, pages 19731–19746. PMLR.
- Jianbin Shen, Junyu Xuan, and Christy Liang. 2023. A determinantal point process based novel sampling method of abstractive text summarization. In 2023 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE.
- Zhao Song, Junze Yin, Lichen Zhang, and Ruizhe Zhang. 2024. Fast dynamic sampling for determinantal point processes. In *International Conference on Artificial Intelligence and Statistics*, pages 244–252. PMLR.
- Zhiliang Tian, Jingyuan Huang, Zejiang He, Zhen Huang, Menglong Lu, Linbo Qiao, Songzhu Mei, Yijie Wang, and Dongsheng Li. 2025. Llm-based rumor detection via influence guided sample selection and game-based perspective analysis. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 28402–28414.
- Huazheng Wang, Jinming Wu, Haifeng Sun, Zixuan Xia, Daixuan Cheng, Jingyu Wang, Qi Qi, and Jianxin Liao. 2024a. Mdr: Model-specific demonstration retrieval at inference time for in-context learning. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4189–4204.
- Liang Wang, Nan Yang, and Furu Wei. 2024b. Learning to retrieve in-context examples for large language models. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1752–1767.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. *Transactions on Machine Learning Research*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

- Zhihua Wen, Zhiliang Tian, Zexin Jian, Zhen Huang, Pei Ke, Yifu Gao, Minlie Huang, and Dongsheng Li. 2024. Perception of knowledge boundary for large language models through semi-open-ended question answering. In *Advances in Neural Information Processing Systems*, volume 37, pages 88906–88931. Curran Associates, Inc.
- Jing Xiong, Zixuan Li, Chuanyang Zheng, Zhijiang Guo, Yichun Yin, Enze Xie, Zhicheng YANG, Qingxing Cao, Haiming Wang, Xiongwei Han, Jing Tang, Chengming Li, and Xiaodan Liang. 2024. DQ-lore: Dual queries with low rank approximation re-ranking for in-context learning. In *The Twelfth International Conference on Learning Representations*.
- Zhao Yang, Yuanzhe Zhang, Dianbo Sui, Cao Liu, Jun Zhao, and Kang Liu. 2023. Representative demonstration selection for in-context learning with two-stage determinantal point process. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5443–5456.
- Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023a. Compositional exemplars for in-context learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 39818–39833. PMLR.
- Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023b. Compositional exemplars for in-context learning. In *International Conference* on Machine Learning, pages 39818–39833. PMLR.
- John M Zelle and Raymond J Mooney. 1996. Learning to parse database queries using inductive logic programming. In *AAAI/IAAI*, pages 1050–1055, Portland, OR. AAAI Press/MIT Press.
- Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active example selection for in-context learning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9134–9148.

A Full Version of Related Work

A.1 Example Selection for ICL

The ICL performance of LLMs depends on the selection of examples. Depending on whether the query information and the task supervision were considered, ICL example selection methods can be divided into three categories: (1) **In-context Insensitive Unsupervised Methods.** These approaches ignore the query information and task supervision. Fu et al. (2022) propose a complexity-based example selection method. Lu et al. (2022) Propose an entropy-based approach to mitigate example order sensitivity. Li and Qiu (2023) use a diversity-guided example search strategy to select examples.

| α | β | GeoQuery | MTOP | NL2Bash | WebQs | RocEnd |
|----------|-------|----------|-------|---------|-------|--------|
| 0.005 | 0.01 | 72.71 | 65.28 | 59.30 | 17.91 | 60.33 |
| 0.005 | 0.05 | 72.36 | 65.14 | 59.76 | 18.11 | 60.33 |
| 0.005 | 0.005 | 72.00 | 65.32 | 59.82 | 18.11 | 59.94 |
| 0.01 | 0.01 | 73.21 | 65.37 | 61.01 | 18.90 | 60.33 |
| 0.01 | 0.05 | 72.36 | 65.19 | 60.07 | 18.21 | 59.94 |
| 0.01 | 0.005 | 72.00 | 65.28 | 59.79 | 18.36 | 59.94 |
| 0.05 | 0.01 | 72.71 | 65.37 | 60.03 | 17.72 | 60.33 |
| 0.05 | 0.05 | 71.29 | 65.28 | 60.03 | 18.06 | 59.94 |
| 0.05 | 0.005 | 71.24 | 65.32 | 58.69 | 18.21 | 59.94 |
| 0.01 | 0 | 70.72 | 65.32 | 59.74 | 17.96 | 59.94 |
| 0 | 0.01 | 69.64 | 65.19 | 59.83 | 17.86 | 59.94 |

Table 6: Experimental results comparing different α and β parameter combinations across multiple datasets.

| | GeoQuery | NL2Bash | MTOP | WebQs | RocEnd |
|---------------|----------|---------|--------|--------|---------|
| dataset size | 404 | 15564 | 7441 | 3778 | 87319 |
| train cost(s) | 61.51 | 863.48 | 370.38 | 155.61 | 7052.34 |
| inf cost(s) | 21 | 180 | 34 | 120 | 1922 |

Table 7: The time cost of our method.

(2) In-context Sensitive Unsupervised Methods.

This category considers query information but ignores the task supervision. Researchers find that selecting different examples can reduce the redundancy of ICL example set (Liu et al., 2022; Agrawal et al., 2023; Hongjin et al., 2022). Wang et al. (2024a) further propose a model-specific example selection method based on feature evaluation to improve ICL performance during inference. Similarly, Liu et al. (2024b) select examples with multiple levels of similarity to queries. (3) In-context Sensitive Supervised Methods. By introducing task supervision, these methods fine-tune ICL example selectors (i.e., retrievers) for more precise example selection. Many studies have improved the quality of ICL examples by iteratively training retrievers (Rubin et al., 2022; Wang et al., 2024b; Li et al., 2023; Liu et al., 2024b). Besides, Xiong et al. (2024) use chain-of-thought generated by LLMs to refine the retriever. Fu et al. (2022) optimize the retriever by calculating semantic similarity, example diversity, and event correlation. To consider diversity, Levy et al. (2023); Yang et al. (2023); Ye et al. (2023b) employ DPP to select diverse example sets. These works only consider relevance to input queries and diversity of examples, our framework further considers relevance among examples.

A.2 Determinantal Point Process (DPP) and Its Applications

Determinantal Point Process (DPP) is a probabilistic model that can select diverse subsets by capturing negative correlations among items of the set.

DPP has seen significant development. Johans-

son et al. (2023) proposed a semi-supervised k-DPP method. Grosse et al. (2024) used a greedy algorithm for k-DPP sampling. To reduce computational complexity, more inference methods were proposed, such as LSMOEA-DPP (Okoth et al., 2022) and Anisotropic DPP (Ghilotti et al., 2024).

DPP is widely used in AI applications, especially for tasks that require diverse sets, such as neural network training (Sheikh et al., 2022), recommendation systems (Liu et al., 2024c), video analysis (Chen et al., 2023), and abstract summary (Shen et al., 2023). DPP also been used to optimize GNN on graph-structured data. (Duan et al., 2022).

Gartrell et al. (2019) propose an extension of DPP called nonsymmetric determinantal point processes (NDPP), which can model both positive and negative correlations among a set of items. Gartrell et al. (2021) reduce NDPP's complexity via kernel decomposition. Han et al. (2022) propose a scalable sampling method for NDPP. Song et al. (2024) propose a fast dynamic algorithm for resampling distributions of NDPP to shorten the sampling time.

B Full Verion of Preliminary

B.1 More about ICL

In-Context Learning (ICL) (Brown et al., 2020) prompts are usually sequences of examples. Given test instance (x_{test}, y_{test}) , LLMs predicts \hat{y} with k-shot ICL prompt:

$$\hat{y} = LLM(e_1 \oplus, ..., \oplus e_k \oplus x_{test})$$
 (11)

Where $e_i = (x_i, y_i)_{i=1}^k$ is the i_{th} example, and \oplus is the concatenation operation. The objective of ICL example selection task is to select k examples

from a pre-constructed example pool such that the predicted value \hat{y} matches its ground truth y_{test} .

B.2 Validity of NDPP's Probability Function

In Eq. 1, the denominator $det(\boldsymbol{L} + \boldsymbol{I}) = \sum_{E \subseteq D} det(\boldsymbol{L}_E)$, i.e., the sum of the determinants of the corresponding sub-kernel matrix of all subsets $E \subseteq D$, must be greater than the numerator $det(\boldsymbol{L}_E)$. According to (Gartrell et al., 2019) Lemma 1, the kernel matrix \boldsymbol{L} is a P0-matrix (all principal minors are nonnegative), which guarantees that the determinant of \boldsymbol{L} and its principal submatrix \boldsymbol{L}_E is nonnegative. And because the principal minors of \boldsymbol{L} are non-negative, the diagonal elements of \boldsymbol{I} are 1, ensuring that the denominator is positive. Thus, Eq. 1 < 1, which is a valid probability value.

B.3 Comparison with DPP on method properties and application scenarios

The kernel matrix of DPP in the traditional setting is restricted to a symmetric positive semi-definite matrix, which can only model the negative correlation between the items in the set, and is more suitable for application scenarios that emphasize the diversity of the subset (e.g., diversity recommendation). NDPP relaxes the symmetry constraint, allowing the kernel matrix to be a nonsymmetric P0-matrix capable of simultaneously modeling both positive and negative correlations, and can be adapted to more complex application scenarios. Experiments on synthetic data in (Gartrell et al., 2019) show that NDPP is more capable of modeling positive and negative correlations between the terms better, whereas DPP would overemphasize the negative correlations between the terms. However, the asymmetry of NDPP may lead to degradation of Fisher information, making training difficult to converge and requiring additional constraint terms.

C Convexity Analysis of the Optimization Objective

Eq. 4 is not concave, because the nonsymmetric kernel matrix results in the Hessian matrix of f in Eq. 4 not being strictly negative definite. Gartrell et al. (2019) shows the Hessian matrix in eq. 8.

D Practical Optimization of the Regularized log-likelihood

We optimize Eq. 8 using the Adam optimizer based on the pseudo-labeled set training set D_{train} constructed in § 4.2, and iteratively optimize the parameter matrices V, B, C until convergence, which is conditional on the rate of change of the log-likelihood of the validation set being less than a preset threshold.

E Approximation Guarantee for Greedy NDPP MAP Inference

In 4.4.2, we present the framework of the greedy algorithm for approximate NDPP MAP inference, which instantiates the classical submodular maximization greedy algorithm (Nemhauser et al., 1978). Gartrell et al. (2021) provided a lower bound on the approximation quality of the greedy algorithm in Theorem 1.

Theorem 1. Consider a nonsymmetric low-rank DPP $L = VV^T + BCB^T$, where V, B are of rank K, and $C \in \mathbb{R}^{K \times K}$. Given a cardinality budget k, let σ_{min} and σ_{max} denote the smallest and largest singular values of L_E for all $E \subseteq D$ and |Y| < 2k. Assume that $\sigma_{min} > 1$. Then,

$$logdet(\mathbf{L}_{E^G}) \ge \frac{4(1 - e^{-1/4})}{2(log\sigma_{max}/log\sigma_{min}) - 1} logdet(\mathbf{L}_{E^*})$$
(12)

where E^G is the output of the greedy algorithm and E^* is the optimal solution of the MAP inference in Eq. 9.

Thus, when the kernel has a small value of $log\sigma_{max}/log\sigma_{min}$, the greedy algorithm finds a near-optimal solution. As mentioned above, there is no evidence that the condition $\sigma_{min}>1$ is usually correct in practice. Gartrell et al. (2021) further provided Corollary 1, which excludes the assumption that $\sigma_{min}>1$ and quantifies this additional term.

Corollary 1. Consider a nonsymmetric low-rank DPP $L = VV^T + BCB^T$, where V, B are of rank K, and $C \in \mathbb{R}^{K \times K}$. Given a cardinality budget k, let σ_{min} and σ_{max} denote the smallest and largest singular values of L_E for all $E \subseteq D$ and |Y| < 2k. Let $\omega := \sigma_{max}/\sigma_{min}$ Then,

$$logdet(\mathbf{L}_{E^{G}}) \ge \frac{4(1 - e^{-1/4})}{2(log\omega)) + 1} logdet(\mathbf{L}_{E^{*}})$$
$$- (1 - \frac{4(1 - e^{-1/4})}{2(log\omega) + 1}) k(1 - log\sigma_{min})$$
(13)

where E^G is the output of the greedy algorithm and E^* is the optimal solution of the MAP inference in Eq. 9. The proof of Theorem 2 and Corollary 1 is given in (Gartrell et al., 2021) appendix F.

F Experimental Setting Details

F.1 Datasets Details

We conduct experiments on 5 text generation tasks, and examples in each dataset are shown in Table 8. We illustrate the details of each dataset as follows.

GeoQuery (Zelle and Mooney, 1996; Shaw et al., 2020) contains a parallel corpus of 880 English questions about US geography paired with Prolog queries. The compositional dataset of GeoQuery was created by Shaw et al. (2020), focusing on compositional generalization.

NL2Bash (Lin et al., 2018) is a dataset for the problem of mapping English sentences to Bash commands. The corpus consists of 9k text-command pairs, where each pair consists of a Bash command scraped from the web and an expert-generated natural language description.

MTOP (Li et al., 2020) is a multilingual parsing dataset with 6 languages. The corpus consists of text-command pairs, where each pair consists of a Bash command scraped from the web and an expert-generated natural language description.

WebQs (Berant et al., 2013) (short for WebQuestions) covers 6,642 question-answer pairs obtained from the web. The questions are selected using the Google Suggest API, and the answers are entities in Freebase.

RocEnd (Mostafazadeh et al., 2016) (short for Roc Ending) is a corpus with 100k stories. The input is a short story consisting of four sentences, and the task objective is to generate the story's ending.

F.2 Metrics Details

Exact Match (EM) (Rajpurkar et al., 2016) is used for GeoQuery, MTOP, and WebQs to measure the accuracy of generated outputs. EM calculates the percentage of predictions that exactly match the ground truth, providing a strict evaluation of correctness.

BLEU-1 (**Papineni et al., 2002**) is applied to RocEnd to assess content alignment in story generation. BLEU-1 focuses on unigram overlap, cap-

turing the overall relevance of generated text to the reference.

BLEU-4 (Papineni et al., 2002) is used for NL2Bash to evaluate structural fidelity in command generation. BLEU-4 emphasizes longer n-gram matches, effectively capturing more complex and syntactically accurate outputs.

F.3 Baselines Details

Random selects non-repeating context examples randomly from the training set, serving as a simple baseline without task-specific guidance.

BM25 (Robertson et al., 2009) retrieves the top-K most similar examples for each test input using the classical sparse retrieval method BM25, which ranks candidates based on low-level textual similarity and selects the highest-scoring ones as context.

EPR (Rubin et al., 2022) leverages a language model to assign positive or negative labels to candidate examples. It then uses the model itself as a scoring function to retrieve effective prompts, selecting the top-K most relevant examples during inference.

CEIL (Ye et al., 2023b) models the probability distribution over the context example subset using Determinantal Point Processes (DPP). It is trained within a contrastive learning framework that balances diversity and relevance through a tunable trade-off parameter, enabling the selection of an optimal example combination.

TTF (Liu et al., 2024b) fine-tunes the retriever using labeled data from the context example set, allowing it to incorporate task-specific modules and better adapt to different tasks through supervised signal.

F.4 Implementation Details

We used GPT-neo-2.7B and GPT-4 as LLM for our study. The maximum context length for the input of the LLM was set at 2048 tokens, and the number of context examples per task was set to 50. If the context size limit of the LLM is exceeded, it will be truncated. We adopted the Adam optimizer with a learning rate of 0.01, and the hyperparameters α and β were both set to 0.01. We perform a grid search using a held-out validation set to select the best-performing hyperparameters. The training was conducted on two NVIDIA A100 GPUs. We initialize the encoder $E_Q(\cdot)$ and $E_Q(\cdot)$ with CEIL

| Dataset | Prompt | Example |
|----------|-------------------|--|
| GeoQuery | {input}\t{output} | Input: what is the population of montana? Output: answer(A,(population(B,A),const(B,stateid(montana)))) |
| Nl2Bash | {input}\t{output} | Input : find all executable files in /home directory. Output : find /home -type f -perm /a=x |
| MTOP | {input}\t{output} | <pre>Input: Create an alarm called 'worktime'. Output: [IN:CREATE_ALARM [SL:ALARM_NAME worktime]]</pre> |
| WebQs | {input} {output} | Input: what does jamaican people speak? Output: Jamaican Creole English Language |
| RocEnd | {input}\t{output} | Input: Dan's parents were overweight. Dan was overweight as well.The doctors told his parents it was unhealthy.His parents understood and decided to make a change.Output: They got themselves and Dan on a diet. |

Table 8: Datasets with corresponding prompts and examples used in the experiments.

| Model | Dataset | GeoQuery | NL2Bash | MTOP | WebQs | RocEnd |
|----------------|-----------------|----------|----------|----------|----------|--------|
| GPT-Neo (2.7B) | Bartlett's Test | 0 | 5.73e-61 | 0 | 7.29e-05 | 0.0251 |
| GPT-4 | Bartlett's Test | 6.92e-03 | 0.0052 | 4.01e-12 | 0.0116 | |

Table 9: The p values of t-test on our method with baselines. The p values are all smaller than 0.05, indicating our improvements are significant.

(Ye et al., 2023a). We employ the implementation from Ye et al. (2023a) for random, BM25, and EPR. For CEIL, we use the result from Liu et al. (2024b) except the result of RocEnd. We also employ the implementation from Ye et al. (2023a) to obtain the result of RocEnd for CEIL.

G Case Study

For an intuitive grasp of the effectiveness of our method, we present some cases on the RocEnd. We compare prompts from Random, BM25, EPR, CEIL, and our method, and show the cases in table 10.

Random method selects examples with disordered topics and irrelevant to the query (e.g., police visits, misplaced blame).

BM25 selects examples based on textual similarity with query and focuses on keyword matching, including redundant content such as the Niagara Falls honeymoon and the Hawaii honeymoon.

EPR also focuses on similarity, but does not ensure emotional consistency. Many examples deviate from the topic (e.g., wedding cancellations). The model learned to generate facts (such as location and time), but did not capture emotions.

CEIL selects more diverse examples, but mixed with negative events such as "luggage arrived 2 days later" and "flight had been cancelled", resulting in completely opposite outcomes.

Compared to baselines:

- 1. Compared to the random method, our method considers relevance to the query, ensuring that all selected examples are related to the query's topic of "travel celebration."
- 2. Compared to BM25 and EPR, our method selects more diverse examples covering a wider range of scenarios (wildlife parks, train museums, city trips, cruise ships), preventing redundant examples while enabling the LLM to capture essential generation patterns beyond keyword matching.
- 3. Compared to CEIL, our method also considers the relevance among examples, selecting examples with uniformly positive outcomes to avoid the negative events, which could mislead the LLM.

Our method considers both the relevance between examples and queries, as well as the relevance and diversity among examples, ensuring that the provided examples maintain relevance to the query while comprehensively covering subsequent story generation patterns and maintaining thematic and emotional consistency, thereby providing appropriate guidance for the LLM to generate the target content.

H Significance Test

We conduct the t-test (Bartlett, 1937) to examine whether the improvements of our method are significant. The p values in Table 9 are all smaller than 0.05, demonstrating the significance of our improvements.

Input

For our 25th wedding anniversary I took my wife on a 2nd honeymoon. We went back to Niagara Falls by train. We stayed on the Canadian side. This is what we did when we first got married.

Example

Random:

- 1.Tom wanted to be healthier. He also loved delicious foods. He decided to eat more fruit. He went to the store.He bought a lot of oranges.
- 2.Dave packed the car. They drove to the resort. They unpacked the car. Something was missing. They blamed Dave.
- 3.It was raining this weekend at Sally's house. She stared out the window for two hours and watched the rainy. Then the police pulled up. The officer banged on the door loudly. Sally answered the door, spoke to the officer, and the officer left.

BM25:

- 1.In August of 2013 my wife and I went to Portland, Maine. We took the train from Boston. We checked into a hotel very close to a baseball stadium. That night we watched the minor league Portland Sea Dogs. Even though they lost, we had a great time in Portland.
- 2.In 1993 we took our kids to Disneyworld. We took the train there. Our first night, my son got sick and threw up. I stayed in the hotel room with him.Luckily he got better the next day and had a great time.
- 3.My brother-in-law took me to Niagara Falls once. We did not bring any ponchos, but we decided to take the ferry. As we got closer to the falls, we got very wet. Everyone on the ferry was laughing at us.Now we know to bring ponchos when we go on the ferry to the falls.

EPR:

- 1.My husband and I went to the San Juan Islands for our honeymoon. We took our small boat and fishing equipment. One morning we went fishing for salmon. We saw orca whales and caught 8 salmon! We were proud of ourselves and had a great day.
- 2.We made plans to go to Nevada. We wanted to visit my husband's mother. She was getting married for the second time. We drove for three days to get there. When we arrived, we were devastated to hear the wedding was cancelled.
- 3.My wife and I are on our honeymoon in Cancun. We have been here for five days and are having a blast. We have been doing nothing but eating, drinking and sleeping. We do miss our two dogs back home though. Thankfully we get to fly back home to see them tomorrow morning.

CEIL

- 1.One of our favorite places to visit as a family is New York City. Once, we only had a day to visit, so we did a lot in a short time. We started out with a ride on the Staten Island Ferry. Then, we made our way all the way up to Central Park.By the time we had to leave, everyone was exhausted.
- 2.Last year this time, we went to NYC to visit our granddaughter. There was a lot of snow in our city that year. Traffic was so bad we missed our train and took a later one. Our city got another storm while we were in NYC.We got back on Sunday and had a hard time getting home.
- 3.We were excited to be flying to Toronto for our first Indian wedding. Our fancy Indian clothes for the events were in our checked luggage. Both of our flights were smooth and landed early. After going through customs we went to collect our luggage. So much for landing early because our luggage arrived 2 days later.

Our:

- 1.Pam took her 2 girls on a train ride. None of them have ever been on a train. The train was headed to Los Angeles. It was a 2 hour train ride.Pam and her girls enjoyed the train ride.
- 2.I took my son to Europe. He said he wanted to go on a train. I agreed to take him on a train. I booked two tickets for a train from Spain to France.My son had a great time on the train!
- 3.I always wanted to go to Saint Croix. For our 25th wedding anniversary, my husband and I went there. We loved the beauty of the land and sea. It is so beautiful there. Now we have decided we are going to move there!

Generated

Random: We went to Niagara Falls by train. **BM25:** We went to Niagara Falls by train.

EPR: We went back to Niagara Falls for our 25th wedding anniversary.

CEIL: We did not want to go back to Niagara Falls.

Our: We had a great time.

Target: We had a great time.

Table 10: Case studies on RocEnd.