

Adaptive Centroid Guided Hashing for Cross-Modal Retrieval

Anonymous submission

Abstract

Deep cross-modal hashing (DCMH) technology is widely used in retrieval tasks due to its low storage costs and efficient computation speed. However, proxy or predefined hash center-based DCMH methods suffer from the following challenges: (1) lack of awareness of data distribution, (2) inability of proxy or predefined hash centers to effectively explore the category semantic information and the unequal scales of each object in a given sample. To address these challenges, we propose a novel adaptive centroid guided hashing (ACGH) method for cross-modal retrieval. Specifically, we extract global and local features using Transformer models, and then fuse them to obtain fine-grained feature representations of multimodal data. Subsequently, the hash centroid generation module leverages the category semantic embedding to construct class hash centers and combine them with learnable label-affinity coefficients (LAC) memory banks to learn adaptive hash centroids. Furthermore, we design a hash centroid guidance module, which uses the hash centroids to guide hash code learning and then updates the hash centers and LAC memory banks through the newly learned hash codes. Extensive experimental results on several benchmark multimodality datasets demonstrate that the proposed ACGH method significantly outperforms other state-of-the-art methods in cross-modal retrieval tasks.

1 Introduction

With the rapid expansion of multimedia data, cross-modal retrieval has become a prominent research area. Its primary goal is to retrieve relevant samples from another modality based on queries from one modality. Cross-modal hashing is an efficient similarity search technology that maps multimodality data into a shared hash space, preserving semantic similarity across heterogeneous data. Moreover, its computational efficiency and reduced storage requirements have made it a crucial technique for large-scale cross-modal retrieval applications.

Deep cross-modal hashing (DCMH) methods can generally be categorized into unsupervised and supervised methods. Unsupervised DCMH [Luo *et al.*, 2021][Mikriukov *et al.*, 2022][Hu *et al.*, 2022] show promising retrieval performance in recent years due to the progress of deep neural networks (DNNs). However, these unsupervised methods suffer from the absence of explicit semantic supervision, which degrades retrieval performance. In contrast, supervised DCMH methods fully utilize label information to capture the semantic relationships between multimodal data more effectively, and thus learn more discriminative hash codes. Most supervision methods [Jiang and Li, 2017] [Xu *et al.*, 2019] [Tu *et al.*, 2022a] [Gu *et al.*, 2019] [Bai *et al.*, 2020] rely on pairwise and triplet strategies to leverage the label information. Specifically, pairwise-based DCMH methods utilize the similarity relationship between data pairs with known labels. The triplet-based method aims to solve local optimization problems by minimizing the distance between anchor points and positive samples, while maximizing the distance between anchor points and negative samples. Recently, several studies have proposed several proxy-based [Huo *et al.*, 2024], [Bai *et al.*, 2023] and hash center-based DCMH methods [Tu *et al.*, 2024] for cross-modal retrieval tasks. These approaches typically use proxies to ensure intra-class compactness and inter-class separation or leverage predefined hash centers to capture the global similarity among multimodal hash codes. However, they still have several limitations in real cross-modal retrieval applications. First, they cannot effectively perceive multimodal data distributions. Second, they fail to fully leverage category semantic information and their semantic relationships.

To overcome these shortcomings, we propose an adaptive centroid guided hashing (ACGH) method for cross-modal retrieval. This method asynchronously learns hash codes and adaptive hash centroids. Specifically, the feature extraction framework of the proposed model is built on the Transformer architecture, thus fully leveraging its powerful feature representation ability to capture more fine-grained features. Meanwhile, we learn adaptive hash centroids by combining the BERTmodel with the LAC memory banks . Additionally, we optimize the LAC memory banks using the newly learned hash codes, achieving forward collaborative optimization between hash centroids and hash codes. Notably, our framework employs an asynchronous learning mechanism to inde-

85 pendently update hash centroids and hash codes, thereby ef-
86 fectively enhancing the model’s overall performance. Exten-
87 sive experimental results demonstrate the superiority of the
88 proposed method in cross-modal retrieval.

89 The contributions of the work can be summarized as fol-
90 lows:

- 91 • We propose a novel DCMH method called Adaptive
92 Centroid Guided Hashing (ACGH). To the best of our
93 knowledge, this is the first work to learn adaptive hash
94 centroids by combining hash centers with their object
95 scales within DCMH methods. These hash centroids effec-
96 tively capture category semantic information and the
97 unequal scales of each label in a given sample.
- 98 • We propose a novel centroid-guided hash code learning
99 strategy, where the hash code learning can be guided by
100 the hash centroids. Then the LAC memory banks and
101 hash centers are updated using the newly generated hash
102 codes. This asynchronous optimization mechanism can
103 learn the optimal hash codes and hash centroids.
- 104 • Extensive experimental results on three benchmark
105 datasets demonstrate the effectiveness of our proposed
106 ACGH approach in cross-modal retrieval tasks.

107 2 Related work

108 2.1 Deep Cross-modal Hashing

109 In recent years, several unsupervised DCMH methods [Liu
110 et al., 2020] [Yu et al., 2021][Tan et al., 2022][Sun et al.,
111 2023][Cui et al., 2024] have achieved excellent performance
112 thanks to the powerful feature extraction ability of DNNs. To
113 fully utilize supervisory information among multimodalities,
114 various supervised DCMH methods [Jiang and Li, 2017][Gu
115 et al., 2019][Shu et al., 2022][Tu et al., 2022a][Qin et al.,
116 2024] [Hu et al., 2024] have been proposed to enhance re-
117 trieval performance. Specifically, DCMHT [Tu et al., 2022a]
118 proposes a differentiable cross-modal hashing method that
119 uses a multi-modal Transformer as the backbone. It gen-
120 erates binary codes through a selection mechanism, allowing
121 for better optimization of the hash codes. DNPH [Qin
122 et al., 2024] employs quadratic spherical mutual informa-
123 tion (QSMI) to minimize neighborhood ambiguity. SCH [Hu et
124 al., 2024] divides sample pairs into three categories based on
125 semantic similarity and then applies different constraints to
126 each category to efficiently utilize the entire hash space. The
127 pairwise or triplet-based strategy is commonly used among
128 these supervised DCMH methods. Therefore, they promote
129 inter-class separation by calculating the relative similarity be-
130 tween samples. However, the intra-class compactness of data
131 across modalities is also weakened, resulting in fuzzy neigh-
132 borhoods.

133 2.2 Proxy or Predefined Center-based Deep 134 Cross-modal Hashing

135 Since several DCMH methods seek to fully explore the global
136 similarity across multimodalities using the proxies and pre-
137 defined hash centers, they generally outperform those based on
138 pairwise and triplet-based approaches. Specifically, DCPH

[Tu et al., 2022b] employs a proxy hashing network to gen-
139 erate proxy hash codes containing category information, and
140 introduces a new edge dynamic softmax loss function that di-
141 rectly uses proxy hash codes as supervisory signals during
142 training. However, classic proxy-based DCMH methods are
143 prone to causing neighborhood ambiguity, leading to subop-
144 timal performance. To solve this issue, PGCH [Bai et al.,
145 2023] constructs a Hadamard matrix to learn the proxies, fur-
146 thermore generating more precise hash codes. TWDH [Tu
147 et al., 2024] designs a novel quantization method to miti-
148 gate information loss during the quantization process. More-
149 over, it proposes a compression approach that better preserves
150 information when reducing feature dimensions, thus effec-
151 tively bridging the gap between high-dimensional and low-
152 dimensional spaces. DNPH [Huo et al., 2024] introduces a
153 uniform distribution constraint, ensuring that each hash bit
154 independently follows a discrete uniform distribution. Al-
155 though the aforementioned methods achieve satisfactory per-
156 formance in some tasks, they are limited by their lack of data
157 distribution awareness and failure to fully leverage category
158 semantics and label relationships.

159 3 Methodology

160 3.1 Notations

161 This work focuses on cross-modal retrieval tasks of image-
162 text. Suppose that we have N training image-text pairs
163 $\mathcal{S} = \{s_i\}_{i=1}^N$, where $s_i = (x_i^v, x_i^t, l_i)$ denotes a image-text
164 pair. Here, $x_i^v \in \mathbb{R}^{d^v}$ and $x_i^t \in \mathbb{R}^{d^t}$ denote the i -th sample
165 from image and text modalities, respectively. $l_i \in \{0, 1\}^{1 \times C}$
166 represents their corresponding multi-label annotation, where
167 C is the category number.

168 3.2 Fine-grained Feature Extraction Module

169 To effectively extract the feature information of multi-
170 modality data, as illustrated in Figure 1, we employ two
171 Transformer encoders, ViT [Dosovitskiy, 2020] and GPT-
172 2 [Radford et al., 2019], to extract the features of the im-
173 age and text modalities, respectively. Specifically, for the i -th
174 image-text pair, the feature representations of the image and
175 text modalities from encoders are given by $f_i^v = \{G_i^v, P_i^v\}$
176 and $f_i^t = \{G_i^t, P_i^t\}$, respectively. Here, $G_i^v \in \mathbb{R}^{1 \times d}$ and $G_i^t \in$
177 $\mathbb{R}^{1 \times d}$ denote global embeddings of image and text modalities,
178 respectively, while $P_i^v \in \mathbb{R}^{T_v \times d}$ and $P_i^t \in \mathbb{R}^{T_t \times d}$ represent
179 their local embedding sequences, respectively, where d is the
180 feature dimension, and T_v and T_t are the number of tokens.
181 The feature extraction process can be represented as follows:
182

$$f_i^{(*)} = E_{\text{trans}}^{(*)}(x_i^{(*)}; \theta^{(*)}), * \in \{v, t\}. \quad (1)$$

183 where $E_{\text{trans}}^{(*)}(\cdot; \cdot)$ is the encoder of Transformer model and
184 $\theta^{(*)}$ represents its parameters. Note that the parameters $\theta^{(*)}$
185 are frozen in this work.

186 After extracting the global and local features of the mul-
187 timodalities using the encoders, we fuse them to obtain the
188 representations of the image and text modalities, respectively.
189 Specifically, we first combine the conceptual token aggrega-
190 tion (CTA) strategy with a single-layer Transformer to cap-
191 ture fine-grained semantic representations from a conceptual
192 perspective. Thus, we have

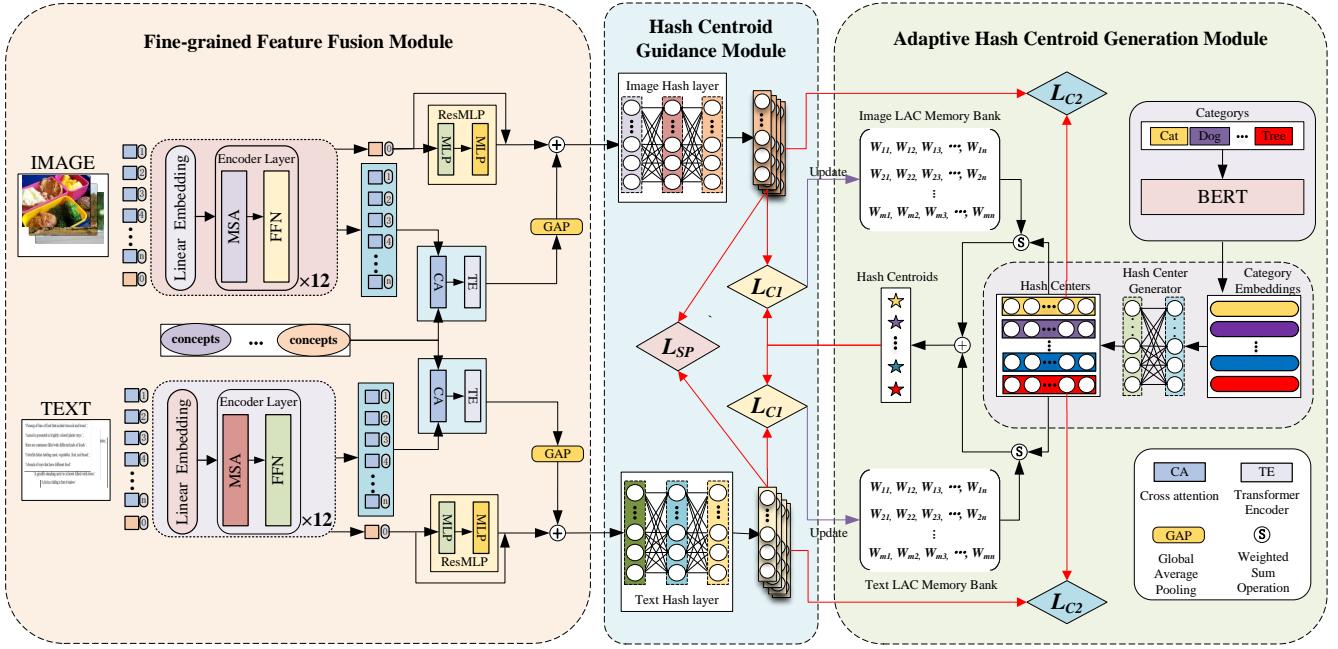


Figure 1: The overall framework of our proposed ACGH approach. It consists of three main modules: the fine-grained feature extraction module, the hash centroid generation module, and the hash centroid guidance module. Specifically, the fine-grained feature extraction module aims to capture intrinsic feature information of multimodalities. The hash centroid generation module combines hash centers with the LAC strategy to generate accurate hash centroids. In the hash centroid guidance module, these centroids guide the learning of hash codes, and the LAC memory bank across modalities is updated using the newly learned hash codes and hash centers.

$$A_i^{(*)} = \text{Softmax} \left(\frac{(QW_q^{(*)})(P_i^{(*)}W_k^{(*)})^T}{\sqrt{d}} \right), \quad (2)$$

$$\tilde{C}_i^{(*)} = A_i^{(*)}(P_i^{(*)}W_v^{(*)}), \quad (3)$$

where $A_i^{(*)} \in \mathbb{R}^{P \times L^{(*)}}$ represents the attention map, and $\tilde{C}_i^{(*)} \in \mathbb{R}^{P \times d}$ is the coarse semantic tokens of multimodalities. $W_q^{(*)}, W_k^{(*)}, W_v^{(*)} \in \mathbb{R}^{d \times d}$ are learnable matrices, and $Q \in \mathbb{R}^{P \times d}$ represents the concept embeddings shared with different modalities.

Afterward, a single-layer Transformer encoder is applied to refine the coarse semantic tokens by capturing their relationships. Therefore, the fine-grained semantic representations of multimodalities are expressed as follows:

$$C_i^{(*)} = E_{st}(\tilde{C}_i^{(*)}; \theta_{est}^{(*)}), \quad (4)$$

where $\theta_{est}^{(*)}$ represents the encoder parameters and $E_{st}(\cdot; \cdot)$ denotes a single-layer Transformer encoder. In this way, CTA transforms variable-length, low-level local embeddings into fine-grained semantic representations of fixed length.

Since global and local information are complementary, we use *ResMLP* to align their dimensions. The specific process is given as follows:

$$r_i^{(*)} = \text{ResMLP}(G_i^{(*)}; \theta_{res}) \in \mathbb{R}^d, \quad (5)$$

where *ResMLP*($\cdot; \cdot$) consists of two identical blocks, each containing one multi-layer perceptron (MLP) with residual

connections. θ_{res} is the trainable weight-shared parameters for image and text modalities.

Then the local and global features with the same dimensions from Eq.(2) and Eq.(5) are fused using an addition operation. Therefore, we can obtain the final representation $X_i^{(*)}$ of the i -th sample as follows:

$$X_i^{(*)} = r_i^{(*)} + \text{GAP}(C_i^{(*)}) \in \mathbb{R}^d. \quad (6)$$

where *GAP*(\cdot) refers to the global average pooling.

3.3 Hash Centroid Guidance Module

The hash centroid guidance module constructs the hashing layer to learn discriminative hash codes. It consists of fully connected layers and a nonlinear activation function, which maps the extracted fused features to continuous hash codes $FU_i^{(*)}$. Thus, we have

$$FU_i^{(*)} = \tanh \left(\text{Hash} \left(X_i^{(*)}; \theta_h^{(*)} \right) \right). \quad (7)$$

where *Hash*(\cdot) represents the hash layer, $\theta_h^{(*)}$ denotes the parameters of the hash layer, and $\tanh(\cdot)$ is the activation function. During the inference phase, we map the continuous hash codes $FU_i^{(*)}$ to binary values of -1 and 1 using the function $H_i^{(*)} = \text{sign}(FU_i^{(*)})$.

3.4 Adaptive Hash Centroid Generation Module

In this module, we first leverage the semantic extraction capability of the BERT [Kenton and Toutanova, 2019] model

232 to embed the category semantic information into a category
 233 embedding matrix $CA = \{ca_i\}_{i=1}^C \in \mathbb{R}^{C \times D}$, where ca_i is
 234 expressed as the embedding vector of the i -th category, and
 235 C represents the category number. D is the dimensionality of
 236 the word embeddings, and v_i is the embedding vector of the
 237 i -th category.

238 Next, we construct a hash center generator, which consists
 239 of two fully connected layers. This generator maps the cat-
 240 egory embedding vectors V into learnable hash centers HC .
 241 This process can be formalized as follows:

$$hc_i = Hash_c(ca_i; \theta_c), \quad (8)$$

$$HC = \{hc_j\}_{j=1}^C \in \mathbb{R}^{C \times K}, \quad (9)$$

242 where θ_c represents the parameters of the hash center gener-
 243 ator $Hash_c(\cdot; \cdot)$, K is the dimensionality of the hash centers,
 244 and hc_i is the hash center for the i -th category.

245 After obtaining the hash centers, we can combine them
 246 with the LAC memory bank to construct the hash centroid
 247 p_i . The calculation process is given as follows:

$$p_i^{(*)} = \sum_{j=1}^C w_{ij}^{(*)} hc_j. \quad (10)$$

248 where $w_{ij}^{(*)}$ denotes the weight of the j -th label for the i -th
 249 sample. Here, $w_i^{(*)} = \{w_{ij}^{(*)}\}_{j=1}^C$ denotes the weight of the i -
 250 th sample in the LAC memory bank, where $\sum_{j=1}^C w_{ij}^{(*)} = 1$.
 251 Note that the initialization of $w_{ij}^{(*)}$ is $w_{ij}^{(*)} = \frac{1}{sum(l_i)}$, where
 252 $sum(l_i)$ is the sum of the one-hot encoded labels for the i -th
 253 sample. The updating process for w_i is given in Section 3.5.

254 3.5 Loss Functions

255 To effectively optimize the hash codes and hash centroids,
 256 we employ an asynchronous learning mechanism that allows
 257 both to be updated independently while adapting to each
 258 other. Once the hash codes $FU_i^{(*)}$ and its corresponding hash
 259 centroid $p_i^{(*)}$ are obtained, we utilize the weighted-softmax
 260 loss L_{WS} to aggregate the samples towards their respective
 261 hash centroids and push them away from irrelevant hash cen-
 262 ters. Thus, we have

$$L_{WS} = -\frac{1}{N} \sum_{i=1}^N \log \frac{ES_P(i)}{ES_N(i)}, \quad (11)$$

$$ES_P(i) = \exp \left(\frac{\cos(FU_i^{(*)}, p_i^{(*)})}{\tau} \right), \quad (12)$$

$$ES_N(i) = \sum_{hc_j \in S_{neg}} \exp \left(\frac{\cos(FU_i^{(*)}, hc_j)}{\tau} \right), \quad (13)$$

263 where S_{neg} represents a subset of hash centers, with each el-
 264 ement corresponding to a hash center hc_j , and j denoting the

index of the 0 element in the one-hot encoded label. τ is a
 265 scale parameter. It is worth noting that this loss is applied in
 266 both the hash code optimization and hash centroid optimiza-
 267 tion processes described below.

268 **Hash Code Learning:** Firstly, to better maintain the pair-
 269 wise semantic similarity between instances, we construct
 270 similarity-preserving loss functions both within and across
 271 modalities, based on the negative log-likelihood loss func-
 272 tion.

273 (1) To effectively preserve the intra-modal similarity, two
 274 asymmetric pairwise negative log-likelihood losses for two
 275 modalities are given as follows:

$$L_{i,j}^{(*)} = -S_{ij} \Omega_{i,j}^{(*)} + \log(1 + e^{\Omega_{i,j}^{(*)}}), \quad (14)$$

$$L_{intra}^{(*)} = \frac{1}{MN} \sum_{i=1}^N \sum_{j=1}^M L_{i,j}^{(*)}, \quad (15)$$

277 where $\Omega_{i,j}^{(*)} = \frac{(FU_i^{(*)})^T (FU_j^{(*)})}{2}$ denotes the inner product of
 278 hash codes within each modality, and S_{ij} is the similarity ma-
 279 trix derived from the known labels.

280 (2) Inter-modal similarity preservation ensures that seman-
 281 tic similarities are maintained across different modalities.
 282 Therefore, we employ the following asymmetric pairwise
 283 negative log-likelihood loss:

$$L_{i,j}^{v2t} = -S_{ij} \Theta_{ij} + \log(1 + e^{\Theta_{ij}}), \quad (16)$$

$$L_{i,j}^{t2v} = -S_{ij} \Phi_{ij} + \log(1 + e^{\Phi_{ij}}), \quad (17)$$

$$L_{inter} = \frac{1}{MN} \sum_{i=1}^N \sum_{j=1}^M (L_{i,j}^{v2t} + L_{i,j}^{t2v}), \quad (18)$$

284 where $\Theta_{ij} = \frac{(FU_i^{(t)})^T (FU_j^{(v)})}{2}$ and $\Phi_{ij} = \frac{(FU_i^{(v)})^T (FU_j^{(t)})}{2}$
 285 indicate the inner products of hash codes across modalities.
 286 Therefore, the overall loss for the semantic similarity perse-
 287 veration of the hash codes can be given as follows:

$$L_{SP} = \beta_1 L_{inter} + \beta_2 L_{intra}^{(*)}, \quad (19)$$

288 where β_1 and β_2 are the weight hyperparameters.

289 Then we use the L_{WS} loss in Eq.(11) to ensure intra-class
 290 compactness and inter-class separation of the hash codes. In
 291 addition, to ensure the quality of the hash codes and prevent
 292 any label weight from dominating the loss, we construct the
 293 quantization loss L_q . Furthermore, the maximum entropy
 294 regularization $R(w)$ is applied to optimize the label weights.
 295 Thus, we have

$$L_q = \frac{1}{N} \sum_{i=1}^N \|H_i^{(*)} - FU_i^{(*)}\|_2^2, \quad (20)$$

$$R(w) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C w_{ij}^{(*)} \log(w_{ij}^{(*)}), \quad (21)$$

296 where $H_i^{(*)}$ denote the binary hash codes corresponding to
 297 the i -th sample and w_{ij} represents the weight of the j -th la-
 298 bel for the i -th sample. Here, the total loss function L_{C1} is
 299 formulated by integrating L_{WS} , L_q , and $R(w)$ as follows:

$$L_{C1} = L_{WS} + \beta_3 L_q + \beta_4 R(w). \quad (22)$$

300 where β_3 and β_4 denote the weight hyperparameters.

301 Finally, the overall objective function L_H for optimizing
 302 the hash codes can be written as follows:

$$L_H = \mu L_{C1} + \lambda L_{SP}, \quad (23)$$

303 where μ and λ are hyperparameters that control the relative
 304 weight of the respective loss terms L_{C1} and L_{SP} .

305 **Adaptive Hash Centroid Learning:** As described above,
 306 the hash centroids are fundamentally derived by combining
 307 the hash centers with the LAC memory banks. Thus, hash
 308 centroid learning consists of two stages: hash code optimiza-
 309 tion and LAC memory bank optimization.

310 **(1) Optimizing hash centers:** To optimize hash centers
 311 and enhance the quality of hash centroids, we first use BERT
 312 to learn category embeddings and then apply a hash center
 313 generator to obtain semantically rich hash centers. To en-
 314 sure that the distribution of the hash center m_{ij} approximates
 315 the distribution of the label word vector s_{ij} , we utilize the
 316 Kullback-Leibler (KL) divergence to construct the label cen-
 317 ter loss L_{CS} . Thus, we have

$$L_{CS} = \sum_{i=1}^C \sum_{j=1}^C s_{ij} \log \frac{s_{ij}}{m_{ij}}, \quad (24)$$

$$s_{ij} = \frac{1}{2} (\cos(ca_i, ca_j) + 1), \quad (25)$$

$$m_{ij} = \frac{1}{2} (\cos(hc_i, hc_j) + 1), \quad (26)$$

318 where ca_i represents the i -th category embedding vector, and
 319 hc_i represents the i -th hash center.

320 Similar to the hash function optimization, we employ L_{WS}
 321 to enhance the similarity (or dissimilarity) of generated hash
 322 centers. Therefore, the total optimization loss for the adaptive
 323 hash centers is formulated as follows:

$$L_{C2} = \alpha_1 L_{WS} + \alpha_2 L_{CS}. \quad (27)$$

324 where α_1 and α_2 denote hyperparameters that control the
 325 contribution of each loss function.

326 **(2) Optimizing LAC memory banks:** To optimize the
 327 LAC memory banks, we compute the gradient of L_{C1} to up-
 328 date the LAC memory bank for both modalities. The opti-
 329 mization process is given as follows:

$$\hat{w}_{ij}^{(*)t} = w_{ij}^{(*)t-1} - \eta \nabla w_{ij}^{(*)}, \nabla w_{ij}^{(*)} = \frac{\partial L_{C1}}{\partial w_{ij}^{(*)t-1}}, \quad (28)$$

330 where $w_{ij}^{(*)t}$ denotes the value of $w_{ij}^{(*)}$ at the t -th iteration.

331 To effectively optimize $w_{ij}^{(*)}$, we employ the Euclidean
 332 projection method [Wang and Carreira-Perpinán, 2013] to
 333 map $\hat{w}_i^{(*)t} = \{w_{ij}^{(*)t}\}_{j=1}^C$ to $w_i^{(*)t} = \{w_{ij}^{(*)t}\}_{j=1}^C$, thereby

334 obtaining the optimal value for each $w_{ij}^{(*)}$. The process is ex-
 335 pressed as follows:

$$w_i^{(*)t} = \min \frac{1}{2} \left\| w_i^{(*)t} - \hat{w}_i^{(*)t} \right\|^2, \sum_{j=1}^C w_{ij}^{(*)} = 1, w_{ij}^{(*)} \geq 0. \quad (29)$$

336 where, $w_i^{(*)t}$ represents the label weight of x_i in the t -th iter-
 337 ation.

338 Finally, after optimizing the hash centers and the LAC
 339 memory banks, we can compute the more precise hash cen-
 340 troid $p_i^t = \frac{1}{2}(p_i^{(v)t} + p_i^{(t)t})$ based on Eq.(10). Therefore, we
 341 adopt an asynchronous mechanism in our model to optimize
 342 both the hash centroids and the hash codes.

4 Experiments

4.1 Experimental Settings

343 **Datasets.** For our experiments, we selected three widely used
 344 datasets for cross-modal hashing retrieval: MIRFLICKR-
 345 25K, MS-COCO, and NUS-WIDE. Their division methods
 346 followed those of previous works.
 347

348 **Baselines.** In this experiment, we conducted a comprehen-
 349 sive comparison with other state-of-the-art DCMH methods,
 350 including both CNN-based and Transformer-based meth-
 351 ods. Specifically, the CNN-based DCMH methods in-
 352 clude AGAH [Gu *et al.*, 2019], PGCH [Bai *et al.*, 2023],
 353 DADH [Bai *et al.*, 2020], and HMAH [Tan *et al.*, 2022]. The
 354 Transformer-based DCMH methods include DCMHT [Tu *et*
 355 *al.*, 2022a], SCH [Hu *et al.*, 2024], DNPh [Qin *et al.*, 2024],
 356 DNPH [Huo *et al.*, 2024], and TWDH [Tu *et al.*, 2024].

357 **Evaluation Metrics.** To better demonstrate the advantages of
 358 our method, we adopted the mean average precision (*mAP*)
 359 of the top 5000 retrieval results to evaluate the retrieval per-
 360 formances of all methods.

4.2 Experimental Results

361 To evaluate the effectiveness of our proposed ACGH method,
 362 we conducted experiments on the MIRFLICKR-25K, MS-
 363 COCO, and NUS-WIDE datasets. We compared our method
 364 with its competitors in two tasks: Image-to-Text (I2T) and
 365 Text-to-Image (T2I). Table 1 shows the detailed compari-
 366 son results of different DCMH methods on three datasets.
 367 It can be seen that our approach achieves the best retrieval
 368 performances across all bit sizes among all methods. Com-
 369 pared to the CNN-based DCMH method, such as AGAH,
 370 DADH, HMAH, and PGCH, the Transformer-based DCMH
 371 approaches achieve significant improvement across differ-
 372 ent bit settings. This is because the feature representation
 373 ability of Transformers is significantly better than that of
 374 CNNs. However, our method still achieves the best retrieval
 375 performance among Transformer-based DCMH approaches.
 376 Specifically, compared to the latest predefined hash center-
 377 based method, TWDH, our method achieves an average im-
 378 provement of 2.03% at 16 bits, 2.94% at 32 bits, and 2.39%
 379 at 64 bits for the I2T task across three datasets. For the T2I
 380 task, the average improvements are 1.88% at 16 bits, 1.97%
 381 at 32 bits, and 1.15% at 64 bits. These results validate that

Methods	Reference	MIRFLICKR-25K			NUS-WIDE			MS-COCO		
		16bits	32bits	64bits	16bits	32bits	64bits	16bits	32bits	64bits
I → T										
AGAH	ICMR'2019	0.7558	0.7815	0.7763	0.7421	0.7821	0.7786	-	-	-
DADH	ICMR'2020	0.8080	0.8280	0.8346	0.6635	0.6700	0.6681	-	-	-
HMAH	TMM'2022	0.8496	0.8771	0.8734	0.7950	0.8002	0.8367	0.6629	0.6693	0.7307
DCMHT	ACMMM'2022	0.8512	0.8715	0.8715	0.7926	0.8021	0.8171	0.7025	0.7477	0.7693
PGCH	TBD'2023	0.7870	0.8151	0.8388	0.7819	0.8190	0.8134	0.5829	0.7389	0.7925
SCH	TPAMI'2024	0.8493	0.8735	0.8913	0.7419	0.7911	0.8113	0.6662	0.7261	0.7754
DNpH	TMM'2024	0.9053	0.9072	0.9111	0.8121	0.8209	0.8246	0.7350	0.7446	0.7454
DNPH	TCSVT'2024	0.8542	0.8788	0.8911	0.7927	0.8121	0.8253	0.7361	0.8089	0.8326
TWDH	TMM'2024	0.8770	0.8682	0.8911	0.8123	0.8287	0.8326	0.7378	0.7992	0.8391
ACGH	Ours	0.9095	0.9195	0.9235	0.8245	0.8450	0.8605	0.7542	0.8199	0.8505
T → I										
AGAH	ICMR'2019	0.7555	0.7664	0.7606	0.7416	0.7672	0.7622	-	-	-
DADH	ICMR'2020	0.7971	0.8169	0.8254	0.6522	0.6547	0.7081	-	-	-
HMAH	TMM'2022	0.8042	0.8248	0.8363	0.7438	0.7724	0.7857	0.7367	0.7780	0.8172
DCMHT	ACMMM'2022	0.8503	0.8647	0.8731	0.7570	0.7608	0.7733	0.6962	0.7478	0.7587
PGCH	TBD'2023	0.7548	0.7900	0.8106	0.7195	0.7762	0.7577	0.5703	0.7595	0.8181
SCH	TPAMI'2024	0.8511	0.8685	0.8802	0.7577	0.7893	0.7996	0.6797	0.7339	0.7770
DNpH	TMM'2024	<u>0.8693</u>	<u>0.8731</u>	0.8759	<u>0.7879</u>	0.7934	0.8010	0.7158	0.7396	0.7489
DNPH	TCSVT'2024	0.8560	0.8682	0.8779	0.7872	<u>0.8097</u>	<u>0.8184</u>	0.7227	<u>0.8159</u>	0.8364
TWDH	TMM'2024	0.8538	0.8638	0.8783	0.7872	0.8039	0.8175	<u>0.7321</u>	0.7980	<u>0.8417</u>
ACGH	Ours	0.8732	0.8812	0.8879	0.8007	0.8258	0.8316	0.7557	0.8170	0.8526

Table 1: Performance comparison of different DCMH methods on three multimodality datasets. The best result is highlighted in bold, and the second-best result is underlined.

the hash centroid guidance strategy, combined with LAC, is more effective than methods relying on predefined hash centers. Furthermore, the top-N precision curve in Figure 2 provides further evidence of the effectiveness of the proposed ACGH method, consistent with the mAP values in Table 1.

4.3 Ablation Study

To assess the contribution of each component in our ACGH framework, we designed four variants to evaluate their impact on the NUS-WIDE dataset. The results of the ablation study are presented in Table 2.

1. ACGH w/o L_{C1} : It was created by removing the loss function L_{c1} . We can observe that its performance is inferior to that of the proposed ACGH method on this dataset. Since it aims to update the LAC memory banks, the experimental results further demonstrate that this loss plays an indispensable role in our model.
2. w/o L_{SP} : This variant was constructed by eliminating the loss function L_{SP} . Its performance is significantly worse than that of our ACGH method. Since L_{SP} is used to preserve the semantic similarity of multimodalities, this suggests that it is crucial for the model's performance.
3. w/o L_{C2} : We constructed this variant of the proposed method by removing the loss function L_{C2} . Similarly, its performance is inferior to that of the ACGH method.

Methods	16bits	32bits	64bits
I → T			
ACGH w/o L_{C1}	0.8227	0.8424	0.8524
ACGH w/o L_{SP}	0.8113	0.8280	0.8359
ACGH w/o L_{C2}	0.8141	0.8373	0.8554
ACGH(PHC)	0.8208	0.8473	0.8377
ACGH	0.8245	0.8450	0.8605
T → I			
ACGH w/o L_{C1}	0.7899	0.8179	0.8260
ACGH w/o L_{SP}	0.7992	0.8138	0.8274
ACGH w/o L_{C2}	0.7950	0.8173	0.8310
ACGH(PHC)	0.8000	0.8086	0.8083
ACGH	0.8007	0.8258	0.8316

Table 2: The ablation study of the proposed method on the NUS-WIDE dataset.

- It is clear that L_{C2} is designed to optimize the hash centers. Therefore, this ablation experiment further validates the effectiveness of continuously updating the hash center in enhancing the retrieval performance.
4. ACGH (PHC): This variant used predefined hash centers (PHC) to replace the adaptive hash centroids. The experimental results show that the proposed ACGH model significantly outperforms the variant in most cases, espe-

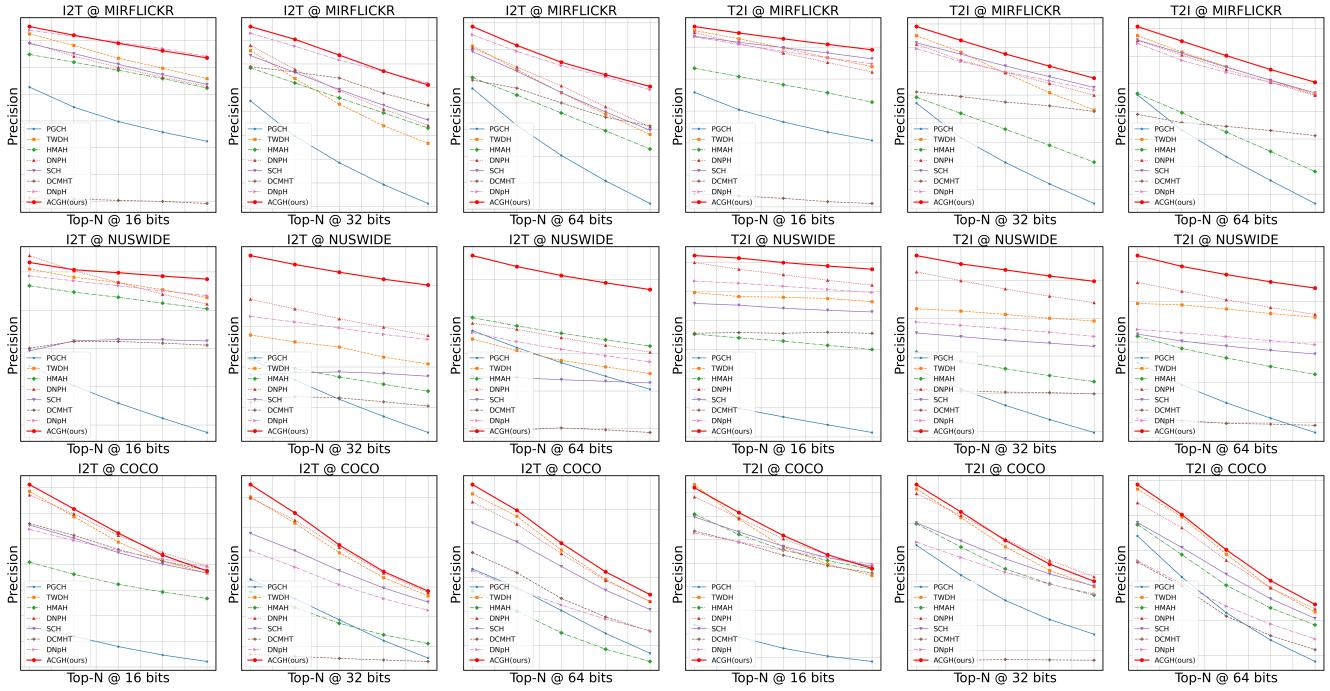


Figure 2: The P@N curves w.r.t. different code lengths on the MIRFLICKR-25K, NUS-WIDE, and MS-COCO datasets.

cially at higher bit sizes. Therefore, it demonstrates that the proposed adaptive hash centroid strategy is superior to the predefined hash center strategy in this framework.

5. ACGH: It can be seen that the full version of our proposed ACGH method outperforms all the aforementioned variants, further demonstrating the effectiveness of each module in our model.

4.4 Parameter Sensitivity

To demonstrate the parameter sensitivity of the proposed method, we conducted experiments on the NUS-WIDE dataset using a 16-bit setting. Specifically, we focused on four parameters in the model: μ , λ , α_1 , and α_2 . Specifically, the parameters μ and λ were set between 0.001 and 10, while the parameters α_1 and α_2 were analyzed within the range of 0.005 to 10. Figure 3 shows the retrieval performance of our ACGH method under different parameter settings. We can see from the results that the proposed ACGH method exhibits excellent stability across a wide range of parameter values on the NUS-WIDE dataset. Therefore, our proposed ACGH method can be easily applied to various retrieval tasks.

5 Conclusion

In this paper, we propose a novel DCMH method termed adaptive centroid guided hashing (ACGH). Initially, we utilize Transformer models to extract global and local features for each modality, which are then fused into a comprehensive fine-grained feature representation. Next, we design two modules: the hash centroid generation module and the hash centroid guidance module. These modules employ an asynchronous optimization mechanism, enabling mutual promo-

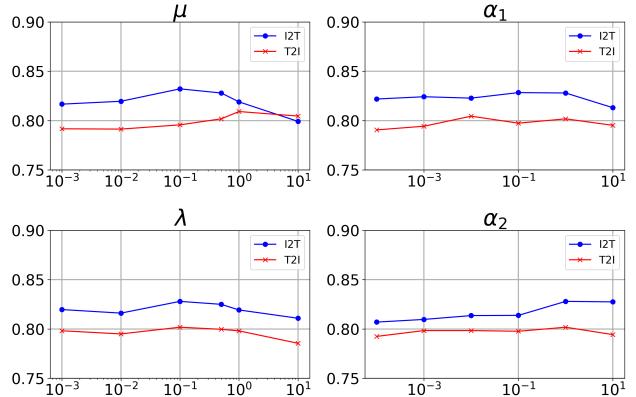


Figure 3: The mAP values of the proposed method under different parameter settings on the NUS-WIDE dataset.

tion and optimization. Specifically, the hash centroid generation module maps the category embeddings obtained from the BERT model to the corresponding hash centers and constructs hash centroids using the LAC strategy. These hash centroids are then used to guide hash code learning. Simultaneously, the newly generated hash codes are used to update the LAC memory banks, thereby improving the accuracy of hash centroids. Therefore, this optimization mechanism enhances the accuracy of both hash codes and hash centroids. Extensive experimental results on three well-known multimodal datasets demonstrate the superiority of the proposed ACGH method in cross-modal retrieval applications.

References

- [Bai *et al.*, 2020] Cong Bai, Chao Zeng, Qing Ma, Jinglin Zhang, and Shengyong Chen. Deep adversarial discrete hashing for cross-modal retrieval. In *Proceedings of the 2020 international conference on multimedia retrieval*, pages 525–531, 2020.
- [Bai *et al.*, 2023] Yibing Bai, Zhenqiu Shu, Jun Yu, Zhengtao Yu, and Xiao-Jun Wu. Proxy-based graph convolutional hashing for cross-modal retrieval. *IEEE Transactions on Big Data*, 2023.
- [Cui *et al.*, 2024] Hui Cui, Lihai Zhao, Fengling Li, Lei Zhu, Xiaohui Han, and Jingjing Li. Effective comparative prototype hashing for unsupervised domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 8329–8337, 2024.
- [Dosovitskiy, 2020] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [Gu *et al.*, 2019] Wen Gu, Xiaoyan Gu, Jingzi Gu, Bo Li, Zhi Xiong, and Weiping Wang. Adversary guided asymmetric hashing for cross-modal retrieval. In *Proceedings of the 2019 on international conference on multimedia retrieval*, pages 159–167, 2019.
- [Hu *et al.*, 2022] Peng Hu, Hongyuan Zhu, Jie Lin, Dezhong Peng, Yin-Ping Zhao, and Xi Peng. Unsupervised contrastive cross-modal hashing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):3877–3889, 2022.
- [Hu *et al.*, 2024] Zhikai Hu, Yiu-ming Cheung, Mengke Li, and Weichao Lan. Cross-modal hashing method with properties of hamming space: A new perspective. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [Huo *et al.*, 2024] Yadong Huo, Qibing Qin, Wenfeng Zhang, Lei Huang, and Jie Nie. Deep hierarchy-aware proxy hashing with self-paced learning for cross-modal retrieval. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [Jiang and Li, 2017] Qing-Yuan Jiang and Wu-Jun Li. Deep cross-modal hashing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3232–3240, 2017.
- [Kenton and Toutanova, 2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacl-HLT*, volume 1, page 2. Minneapolis, Minnesota, 2019.
- [Liu *et al.*, 2020] Song Liu, Shengsheng Qian, Yang Guan, Jiawei Zhan, and Long Ying. Joint-modal distribution-based similarity hashing for large-scale unsupervised deep cross-modal retrieval. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, pages 1379–1388, 2020.
- [Luo *et al.*, 2021] Xiao Luo, Daqing Wu, Zeyu Ma, Chong Chen, Minghua Deng, Jianqiang Huang, and Xian-Sheng Hua. A statistical approach to mining semantic similarity for deep unsupervised hashing. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 4306–4314, 2021.
- [Mikriukov *et al.*, 2022] Georgii Mikriukov, Mahdyar Ravanbakhsh, and Begüm Demir. Unsupervised contrastive hashing for cross-modal retrieval in remote sensing. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4463–4467. IEEE, 2022.
- [Qin *et al.*, 2024] Qibing Qin, Yadong Huo, Lei Huang, Jiangyan Dai, Huihui Zhang, and Wenfeng Zhang. Deep neighborhood-preserving hashing with quadratic spherical mutual information for cross-modal retrieval. *IEEE Transactions on Multimedia*, 2024.
- [Radford *et al.*, 2019] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [Shu *et al.*, 2022] Zhenqiu Shu, Yibing Bai, Donglin Zhang, Jun Yu, Zhengtao Yu, and Xiao-Jun Wu. Specific class center guided deep hashing for cross-modal retrieval. *Information sciences*, 609:304–318, 2022.
- [Sun *et al.*, 2023] Lina Sun, Yewen Li, and Yumin Dong. Learning from expert: Vision-language knowledge distillation for unsupervised cross-modal hashing retrieval. In *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*, pages 499–507, 2023.
- [Tan *et al.*, 2022] Wentao Tan, Lei Zhu, Jingjing Li, Huaxiang Zhang, and Junwei Han. Teacher-student learning: Efficient hierarchical message aggregation hashing for cross-modal retrieval. *IEEE Transactions on Multimedia*, 25:4520–4532, 2022.
- [Tu *et al.*, 2022a] Junfeng Tu, Xueliang Liu, Zongxiang Lin, Richang Hong, and Meng Wang. Differentiable cross-modal hashing via multimodal transformers. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 453–461, 2022.
- [Tu *et al.*, 2022b] Rong-Cheng Tu, Xian-Ling Mao, Rong-Xin Tu, Binbin Bian, Chengfei Cai, Hongfa Wang, Wei Wei, and Heyan Huang. Deep cross-modal proxy hashing. *IEEE Transactions on Knowledge and Data Engineering*, 35(7):6798–6810, 2022.
- [Tu *et al.*, 2024] Junfeng Tu, Xueliang Liu, Yanbin Hao, Richang Hong, and Meng Wang. Two-step discrete hashing for cross-modal retrieval. *IEEE Transactions on Multimedia*, 2024.
- [Wang and Carreira-Perpiñán, 2013] Weiran Wang and Miguel A Carreira-Perpiñán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. *arXiv preprint arXiv:1309.1541*, 2013.
- [Xu *et al.*, 2019] Ruiqing Xu, Chao Li, Junchi Yan, Cheng Deng, and Xianglong Liu. Graph convolutional network hashing for cross-modal retrieval. In *Ijcai*, volume 2019, pages 982–988, 2019.

- 568 [Yu *et al.*, 2021] Jun Yu, Hao Zhou, Yibing Zhan, and
569 Dacheng Tao. Deep graph-neighbor coherence preserving
570 network for unsupervised cross-modal hashing. In *Pro-*
571 *ceedings of the AAAI conference on artificial intelligence*,
572 volume 35, pages 4626–4634, 2021.