

基于深度学习的1-to-K的跨模态跨语言信息检索引擎

一、背景与意义

随着信息化进程的加速，数据量呈指数级增长，不同模态（如文本、图像、视频、音频等）的数据不断涌现。这些数据往往涉及多种语言和多种模态，使得用户在获取信息时面临复杂的跨模态和跨语言障碍。例如，一个用户可能希望通过图像检索相关的文本信息，或者通过母语搜索外语的多媒体内容。传统的单模态检索（仅限于文本或图像）已无法满足这些需求，因此，跨模态跨语言信息检索系统应运而生。

跨模态跨语言信息检索的目标是打破模态和语言之间的障碍，允许用户以一种模态和语言输入查询，并在不同模态和语言的数据集上找到对应的匹配结果。例如，用户可以使用中文文本搜索英语图片描述，或通过图像查询与其内容相符的多语言文本。这样的系统在多媒体搜索、跨语言知识问答、跨文化电商推荐等应用场景中具有广泛的应用前景。

二、现状分析

目前，跨模态跨语言检索技术主要面临以下挑战：

- 语义鸿沟:** 不同模态数据之间存在着天然的语义鸿沟，如何有效地将不同模态的特征进行对齐，是跨模态检索的核心问题。
- 数据稀疏性:** 多模态数据往往存在数据稀疏性问题，难以通过传统的机器学习方法进行有效的特征学习。
- 样本增量学习:** 在检索数据库中,一旦加入新数据,则需要花费大量时间重新训练模型或者重新计算,此时样本的增量学习就显得尤为重要。

现有的跨模态检索技术主要分为以下几类：

- 传统方法:** 使用手工设计的特征提取技术，将不同模态的数据映射到一个共同的空间。适合小规模数据，效率高但效果有限。
- 基于深度学习的方法:** 利用神经网络自动提取特征，实现图像和文本的跨模态匹配。适合大规模数据，精度高但需要大量标注数据。
- 基于哈希编码的方法:** 将数据映射到二进制码空间，实现快速检索。适用于大数据场景，存储和检索效率高，但精度略低。

三、方案设计

设计需求

• 技术可行性

该跨模态跨语言信息检索系统的核心需求是提供高效、精准的跨模态查询服务，即用户可通过文本或图像输入，实现跨模态内容的推荐与检索。系统需具备在大规模数据上提取语义关联、捕捉用户偏好并适应不同模态不同语言数据的能力。

现代深度学习技术的发展为此类系统的实现提供了稳固的技术基础。文本和图像嵌入模型、异构图神经网络、注意力机制等技术在大规模、多模态数据处理中的表现已在各类应用中得到广泛验证，并形成较为成熟的框架支持。此外，通过优化深度学习的计算效率和资源管理，使得模型训练和系统运行可在企业级硬件环境下实现。

• 经济可行性

此跨模态跨语言信息检索系统能够显著提升检索效率，减少用户在大量内容中手动筛选的时间，同时满足对用户偏好的精准理解，从而提高平台的用户满意度与活跃度，促进流量转化。该系统在数据管理方面对开发团队提供便捷的操作界面，使后续维护成本较低。尽管系统开发初期需大量投入，尤其是深度学习模型训练和大规模数据处理成本，但对于依赖信息检索的大型平台而言，开发成本是可承受的，并且系统产生的收益将高于其开发及维护费用。因此经济上是可行的。

• 操作可行性

跨模态跨语言信息检索系统依托现代计算设备和网络技术，用户仅需输入文本或上传图像，即可获得相关内容的推荐。操作便捷、界面友好，使得用户无需复杂的学习过程即可掌握。此外，系统支持后台数据的自动更新和模型定期优化，管理员仅需定期维护数据库即可。因此在操作上具有良好的可行性。

产品功能

此跨模态跨语言信息检索系统具备以下核心功能：

- 跨模态跨语言查询支持：**用户可通过输入文本或上传图像实现跨模态检索，支持用户输入不同语言的文本内容，系统将自动识别并翻译成标准语言进行检索，查询结果将显示与用户输入相关的文本和图像，确保用户能够轻松找到所需的信息。
- 个性化推荐：**基于用户历史交互行为与当前检索内容，系统结合用户长期和短期偏好，提供个性化的跨模态推荐，提升推荐的精准度与用户体验。
- 语义关联挖掘：**系统自动识别和计算文本与图像之间的深层语义关联，确保跨模态推荐结果在语义上与用户输入紧密相关，提升检索的准确性。
- 注意力加权机制：**通过对文本和图像特征加权分析，系统可动态调整不同模态信息在检索中的权重，从而更好地匹配用户检索意图。
- 数据稀疏问题缓解：**采用先进的对比学习机制，提升系统在数据稀疏情况下的泛化能力，确保即使在数据量不足时也能生成可靠的推荐结果。
- 多任务优化：**通过多任务优化策略解决模型训练中的冲突问题，实现多个目标任务的协同训练，从而在不同模态任务下均表现优异，提高整体检索效果。

应用场景

该跨模态跨语言检索系统可以用在下面的应用场景中。

- 电子商务推荐系统：用户可以通过上传商品图片或输入描述快速找到相关商品。
- 内容审核与监控：社交媒体平台可以利用该系统对用户上传的文本和图片内容进行合规性审核，减少人工审核压力，保证审核的高效性和一致性。
- 文化遗产保护：文物研究人员可以借助该系统通过图片和文本快速检索文物信息，从而便于文物的研究与记录，促进文化保护。
- 翻译与本地化：对于多媒体内容的本地化（如视频字幕和网站内容），跨模态跨语言系统可以识别不同语言的图文内容并提供精准翻译，从而帮助企业更好地向全球市场推广产品或内容。
- 新闻聚合与实时监控：针对国际新闻和突发事件，系统可以从不同语言的文本和图片中提取信息，为用户提供全面的新闻汇总和多模态分析。这在媒体和政府部门的国际舆情监控中具有重要应用价值。

四、技术路线

1. 跨模态嵌入建模与1-to-K对比学习

- **文本嵌入**：使用 **XLM-R (Cross-lingual Language Model by RoBERTa)** 作为多语言文本编码器。XLM-R 是一种基于 Transformer 的多语言预训练模型，能够处理 100 多种语言。它在大规模多语言数据集上进行了预训练，能有效学习不同语言之间的语义一致性。
 - 输入：
 - 源语言句子：以源语言的方式表示的文本句子，如用原始文本或子词（subword）表示的语言句子。
 - 目标语言句子：以目标语言的方式表示的文本句子，如用原始文本或子词表示的语言句子。
 - 输出：
 - 源语言句子的文本表示：由XLM-R模型生成的源语言句子的固定长度向量表示。
 - 目标语言句子的文本表示：由XLM-R模型生成的目标语言句子的固定长度向量表示。
 - 假设输入文本序列为 $T = \{t_1, t_2, ..., t_n\}$ ，XLM-R将其映射为固定长度的语义向量 \mathbf{v}_T ，用于跨语言的语义对齐，即：

$$\mathbf{v}_T = \text{XLM-R}(T)$$

- **图像嵌入**：使用 **Swin Transformer** 作为图像编码器。Swin Transformer 是一种基于滑动窗口机制的视觉 Transformer 模型，它在多个计算机视觉任务中表现优异，能够提取出具有良好语义表示的图像特征。Swin Transformer 将输入图像分割成不重叠的图块（patch），并通过滑动窗口的自注意力机制逐层提取特征，最终得到图像的全局特征向量表示。该特征向量在跨模态对比学习中与文本向量对齐，从而在同一语义空间内进行跨模态检索。

- 对于输入图像 I ，经过图像嵌入模型映射为高维向量 \mathbf{v}_I ：

$$\mathbf{v}_I = \text{Swin-Transformer}(I)$$

将文本向量 \mathbf{v}_T 和图像向量 \mathbf{v}_I 投射到统一的语义空间，便于跨模态检索。

- 1-to-K对比学习：**在跨模态检索中，单一语言对齐可能引发误差传播和优化方向偏差。因此，引入1-to-K对比学习方法，将每个图像嵌入同时与多个语言的文本嵌入对齐，从而消除以往使用单一语言（如英语）进行对齐所带来的问题。在1-to-K对比学习中，图像嵌入 \mathbf{v}_I 同时与多语言的文本嵌入 $\{\mathbf{v}_T^1, \mathbf{v}_T^2, \mathbf{v}_T^3, \dots, \mathbf{v}_T^K\}$ 对齐，消除单一语言对齐导致的误差传播和方向偏差。

对比损失函数为：

$$\mathcal{L}_{\text{对比}} = -\log \frac{\sum_{k=1}^K \exp(\text{simcos}(\mathbf{v}_I, \mathbf{v}_T^k))}{\sum_{k=1}^K \exp(\text{simcos}(\mathbf{v}_I, \mathbf{v}_T^k)) + \sum_j \exp(\text{simcos}(\mathbf{v}_I, \mathbf{v}_{T_j}^-))}$$

这里的simcos表示余弦相似度，定义如下：

$$\text{simcos}(\mathbf{v}_T, \mathbf{v}_I) = \frac{\mathbf{v}_T \cdot \mathbf{v}_I}{\|\mathbf{v}_T\| \|\mathbf{v}_I\|}$$

通过1-to-K对比学习，模型能够在同一语义空间内将图像与多种语言的文本对齐，从而提升跨语言检索的一致性和鲁棒性。

3. 跨模态注意力机制

为了有效捕捉跨模态信息的相关性，引入跨模态注意力机制来对图像和文本特征进行加权融合。注意力机制可以为不同语言的文本特征分配不同的权重，从而在模态间实现精确的语义组合。

对于每个图像嵌入 \mathbf{v}_I 和多语言文本嵌入 \mathbf{v}_T^k ，使用注意力机制为每个文本嵌入计算注意力权重 α_k ，以控制其在融合表示中的贡献：

$$\alpha_k = \frac{\exp(f(\mathbf{v}_T^k, \mathbf{v}_I))}{\sum_{j=1}^K \exp(f(\mathbf{v}_T^j, \mathbf{v}_I))}$$

这里 $f(\cdot)$ 为评分函数，通常使用点积或前馈神经网络来量化文本和图像嵌入之间的相似度。最终的模态融合表示为：

$$\mathbf{v}_{\text{融合}} = \sum_{k=1}^K \alpha_k \cdot \mathbf{v}_T^k + \beta \cdot \mathbf{v}_I$$

这里， β 是图像特征的全局权重。该加权机制有助于提升不同模态特征间的互动，以更好地捕捉跨模态语义一致性。

4. 语义相似度计算

在统一语义空间中使用余弦相似度或欧氏距离来计算文本和图像的语义相似性，以便检索匹配。对于文本嵌入 \mathbf{v}_T 和图像嵌入 \mathbf{v}_I ，余弦相似度 simcos 和欧氏距离 d 分别定义为：

$$\text{simcos}(\mathbf{v}_T, \mathbf{v}_I) = \frac{\mathbf{v}_T \cdot \mathbf{v}_I}{\|\mathbf{v}_T\| \|\mathbf{v}_I\|}$$

$$d(\mathbf{v}_T, \mathbf{v}_I) = \|\mathbf{v}_T - \mathbf{v}_I\|$$

5. 对比学习缓解数据稀疏问题

采用星形对比学习，将同类模态数据映射到一个紧凑的语义空间，以提升模型在稀疏数据下的泛化能力。给定正样本对 $(\mathbf{v}_T^+, \mathbf{v}_I^+)$ 和负样本对 $(\mathbf{v}_T^-, \mathbf{v}_I^-)$ ，定义对比学习损失 $\mathcal{L}_{\text{对比}}$ ：

$$\mathcal{L}_{\text{对比}} = -\log \frac{\exp(\text{simcos}(\mathbf{v}_T^+, \mathbf{v}_I^+))}{\exp(\text{simcos}(\mathbf{v}_T^+, \mathbf{v}_I^+)) + \sum_k \exp(\text{simcos}(\mathbf{v}_T, \mathbf{v}_{I_k}^-))}$$

6. 多任务学习与梯度冲突优化

在多任务学习中，利用梯度手术技术调整不同任务的梯度，避免梯度冲突。设任务 A 和任务 B 的梯度分别为 $\nabla \mathcal{L}_A$ 和 $\nabla \mathcal{L}_B$ ，当梯度内积 $\nabla \mathcal{L}_A \cdot \nabla \mathcal{L}_B < 0$ 时，进行梯度裁剪，以减小梯度冲突：

$$\nabla \mathcal{L}_{\text{优化}} = \nabla \mathcal{L}_A + \text{proj}_{\nabla \mathcal{L}_A}(\nabla \mathcal{L}_B)$$

其中 $\text{proj}_{\nabla \mathcal{L}_A}(\nabla \mathcal{L}_B)$ 为 $\nabla \mathcal{L}_B$ 在 $\nabla \mathcal{L}_A$ 上的投影，通过该方法，确保不同任务间的梯度不相互干扰，从而提升多任务学习的有效性。

7. 模型优化与实时检索服务

- 离线训练：**基于大规模数据集进行模型离线训练，通过上述技术路径不断优化文本和图像的跨模态嵌入表示。
- 在线推理：**将训练好的模型部署至在线服务端，实时处理用户的文本或图像输入，并利用前述语义相似度匹配机制返回最优检索结果，提供高效的跨模态检索体验。

上述技术路线确保了系统对不同模态的兼容性与扩展性，提升跨模态检索的准确性和实时响应性能。

五、总结

本设计文档提出了一种基于深度学习的跨模态跨语言检索系统，该系统能够有效地解决不同模态不同语言数据之间的语义鸿沟问题，从而更好地满足用户的信息需求。未来，我们可以进一步探索以下方向：

- **探索更有效的融合方法:** 例如动态融合或对抗学习, 进一步提升融合效果。
- **构建大规模的多模态数据集:** 为模型训练提供更丰富的数据资源。
- **开发更友好的用户界面:** 提升用户体验。

六、参考文献

1. 王宏志, 燕钰. 深度学习驱动的跨模态数据检索. *Journal of Harbin University of Science & Technology*. 2021;26(1):9-16. doi:10.15938/j.jhust.2021.01.002
2. FENG Xia, HU Zhi-yi, LIU Cai-hua. [Survey of Research Progress on Cross-modal Retrieval](#)[J]. *Computer Science*, 2021, 48(8): 13-23. <https://doi.org/10.11896/jsjcx.200800165>
3. Zhijie Nie, Richong Zhang, Zhangchi Feng, Hailang Huang, and Xudong Liu. 2024. Improving the Consistency in Cross-Lingual Cross-Modal Retrieval with 1-to-K Contrastive Learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*. Association for Computing Machinery, New York, NY, USA, 2272–2283. <https://doi.org/10.1145/3637528.3671787>
4. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2020). Unsupervised Cross-lingual Representation Learning at Scale. arXiv preprint arXiv:1911.02116, cs.CL. Retrieved from <https://arxiv.org/abs/1911.02116>
5. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., & Guo, B. (2021). Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. arXiv preprint arXiv:2103.14030, cs.CV. Retrieved from <https://arxiv.org/abs/2103.14030>