

# 画像検索精度向上のための最新研究動向レビュー

## 1. プロンプト誘導型注意ヘッド選択 (PHS)

論文: "Prompt-Guided Attention Head Selection for Focus-Oriented Image Retrieval" (2025)

**主要アイデア:** - ユーザーの視覚的プロンプト（点、ボックス、セグメンテーション）に基づいて、Vision Transformer (ViT)の特定の注意ヘッドを選択 - 複雑な画像から特定のオブジェクトに焦点を当てた検索を実現 (Focus-Oriented Image Retrieval: FOIR) - モデルの再学習不要、入力画像の改変も不要

**技術的特徴:** - ViTの最終層の注意ヘッドが持つ高レベルな顕著性・セグメンテーション情報を活用 - 各注意ヘッドの注意マップとユーザープロンプトのマッチングによる選択 - 人間とモデルの視覚的理解のギャップを埋める高レベル知覚マッチング機構

**MSBAへの応用可能性:** - MSBAのバッチ内注意機構に選択的注意の概念を導入可能 - バッチ内の特定オブジェクトに焦点を当てた特徴抽出が可能に - ユーザー指定の関心領域に基づく動的な注意ヘッド選択機構の統合

## 2. デュアルブランチTransformer (CT-Tran)

論文: "Complementary two-branch Transformer for multi-label image retrieval" (2025)

**主要アイデア:** - 視覚特徴とラベル間の複雑な依存関係を自律的に発見するデュアルブランチアーキテクチャ - Vision Transformerブランチとラベル埋め込みブランチの相補的設計 - ランダムマスキング戦略と組み合わせたマルチヘッド自己注意機構

**技術的特徴:** - CNNバックボーンを使わない純粋なTransformerベースのモデル - クロスエントロピー損失とトリプレット損失の組み合わせによる識別的特徴学習 - マルチラベル画像の複雑なセマンティクス関係のモデリング

**MSBAへの応用可能性:** - MSBAにデュアルブランチ構造を導入し、バッチ内画像とラベル情報の相互作用を強化 - ランダムマスキング戦略によるバッチ内情報の堅牢な学習 - マルチラベル依存関係のモデリング手法をバッチ内画像間の関係性モデリングに応用

### 3. マルチモーダルLLMによるスパース表現

**論文:** "Rethinking Sparse Lexical Representations for Image Retrieval in the Age of Rising Multi-Modal Large Language Models" (2024)

**主要アイデア:** - マルチモーダル大規模言語モデル（M-LLM）を活用した画像特徴の抽出とテキストデータへの変換 - 自然言語処理で使用される効率的なスパース検索アルゴリズムの画像検索への応用 - キーワードベースの画像検索における反復的なクエリ改善

**技術的特徴:** - データ拡張技術によるキー拡張とその効果の分析 - 画像とテキストデータ間の関連性を評価するメトリクスの導入 - 従来のビジョン言語モデルベースの手法と比較した優れた精度とリコール性能

**MSBAへの応用可能性:** - MSBAにスパース表現の概念を導入し、計算効率と検索精度のバランスを最適化 - M-LLMから抽出した意味情報をバッチ内注意機構に統合 - キーワード拡張技術をバッチ内コンテキスト強化に応用

### 4. クロスウィンドウ自己注意（CSWin）

**論文:** "Vision Transformers (ViT) in Image Recognition: Full Guide" (2024)

**主要アイデア:** - 画像の異なる部分を同時に分析するクロス形状のウィンドウ自己注意機構 - 従来の自己注意機構と比較して大幅な処理速度向上 - 局所的・大域的特徴の効率的な統合

**技術的特徴:** - 水平・垂直方向のクロス形状ウィンドウによる注意計算 - 計算複雑性の削減と並列処理の最適化 - 異なるスケールの特徴間の効率的な情報伝播

**MSBAへの応用可能性:** - MSBAのマルチスケールスライディングウィンドウモジュールにクロスウィンドウ構造を導入 - バッチ内の画像間の水平・垂直方向の関係性モデリングを強化 - 計算効率と注意範囲のバランスを最適化

### 5. 因果注意Transformer

**論文:** "Causal Attention Transformer for Video Text Retrieval" (2025)

**主要アイデア:** - 因果推論ネットワークによるビデオテキストペア間の因果特徴抽出 - 時間的依存関係と意味的関連性の両方をモデル化 - 偽相関の影響を軽減し、本質的な特徴関係に焦点

**技術的特徴:** - 因果グラフ構造に基づく注意機構の設計 - 介入と反実仮想に基づく特徴学習 - 時間的一貫性と意味的整合性の同時最適化

**MSBAへの応用可能性:** - バッチ内画像間の因果関係モデリングによる擬似相関の排除 - 双方向コンテキスト伝播モジュールに因果推論機構を統合 - バッチ内の本質的な特徴関係に焦点を当てた学習の強化

## 6. 総合的な技術トレンドと応用可能性

1. **選択的注意機構:**
2. ユーザープロンプトや画像コンテンツに基づく動的な注意ヘッド選択
3. 特定のオブジェクトや領域に焦点を当てた特徴抽出
4. MSBAへの応用: バッチ内の関連画像に選択的に注意を向ける機構
5. **マルチモーダル情報統合:**
6. 視覚特徴とテキスト/ラベル情報の相補的な活用
7. 異なるモダリティ間の複雑な依存関係のモデリング
8. MSBAへの応用: バッチ内画像の意味的関係性の強化
9. **効率的な注意計算:**
10. クロスウィンドウ構造やスパース表現による計算効率の向上
11. 局所的・大域的特徴の効率的な統合
12. MSBAへの応用: バッチサイズ拡大と計算コスト削減の両立
13. **因果推論と本質的特徴抽出:**
14. 偽相関の排除と本質的な特徴関係への焦点
15. 介入と反実仮想に基づく堅牢な特徴学習
16. MSBAへの応用: バッチ内の本質的な類似性に基づく特徴学習
17. **反復的クエリ改善:**
18. ユーザーフィードバックや初期検索結果に基づくクエリの反復的改善
19. 検索意図の明確化と検索精度の向上
20. MSBAへの応用: バッチ内情報に基づく特徴表現の反復的改善