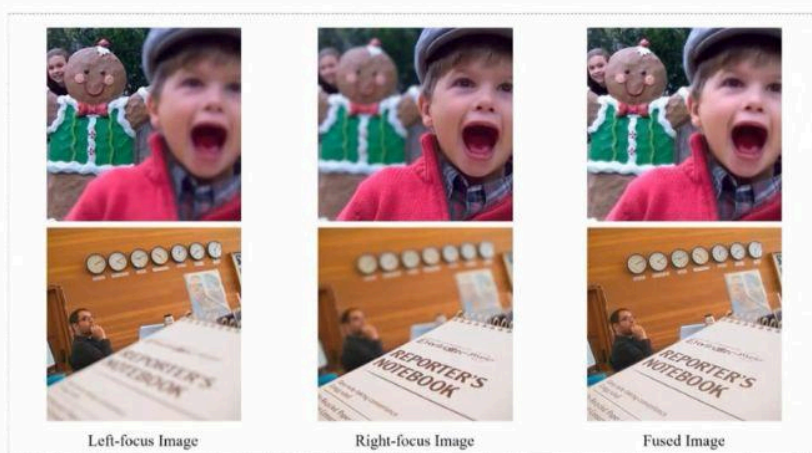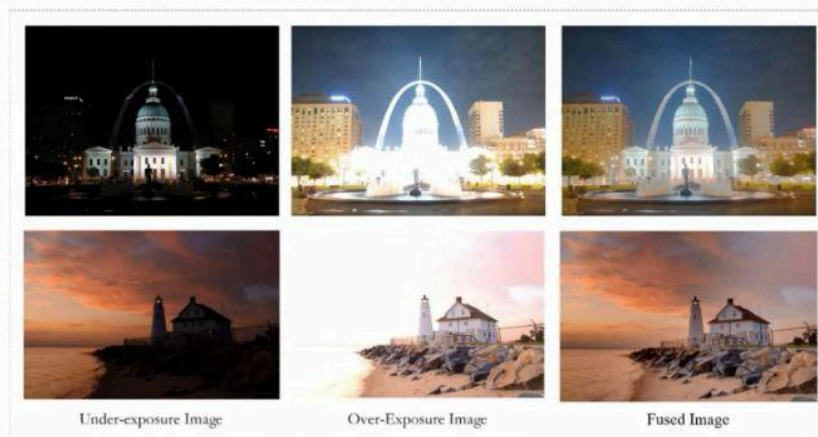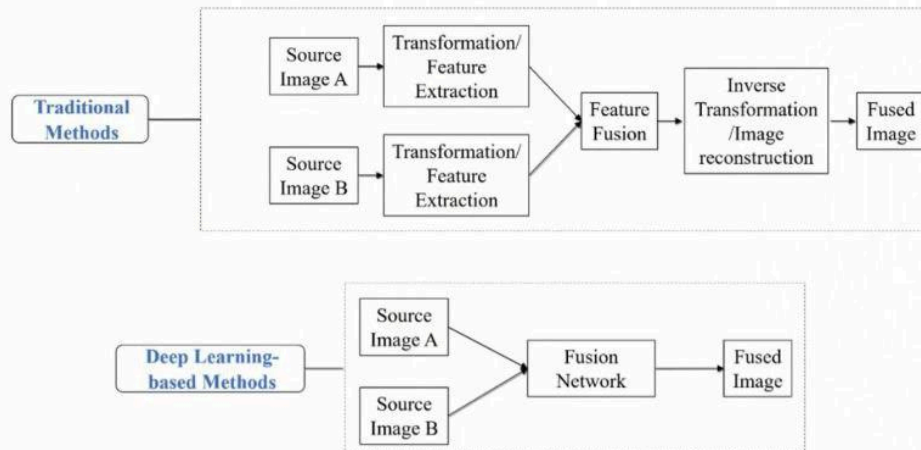# Image Fusion

- **Definition**: A technique where multiple images are combined into a single image to enhance overall information or quality.
- **Purpose**: To create a composite image with more detailed or comprehensive information than any single input image.
- **Types of Image Fusion**:
  - **Pixel-level fusion**:
    - Merges pixel values from different images.
    - Techniques include averaging, PCA, and wavelet transforms.
  - **Feature-level fusion**:
    - Combines extracted features (e.g., edges, textures, shapes).
    - Focuses on combining relevant features from the images.
  - **Decision-level fusion**:
    - Merges decisions made on each image (e.g., classifications or object detections).
- **Goal**:
  - Improve clarity, information content, or visualization in the combined image.
  - Useful when working with images taken under different conditions or sensors.
- **RGB-IR Image Fusion**:
  - process of combining **RGB** images with **Infrared (IR)** images to create a new image that enhances both visible and infrared information.
  - This fusion technique leverages the strengths of both image types for applications where both visible and thermal details are required.

# Multi-exposure Image Fusion
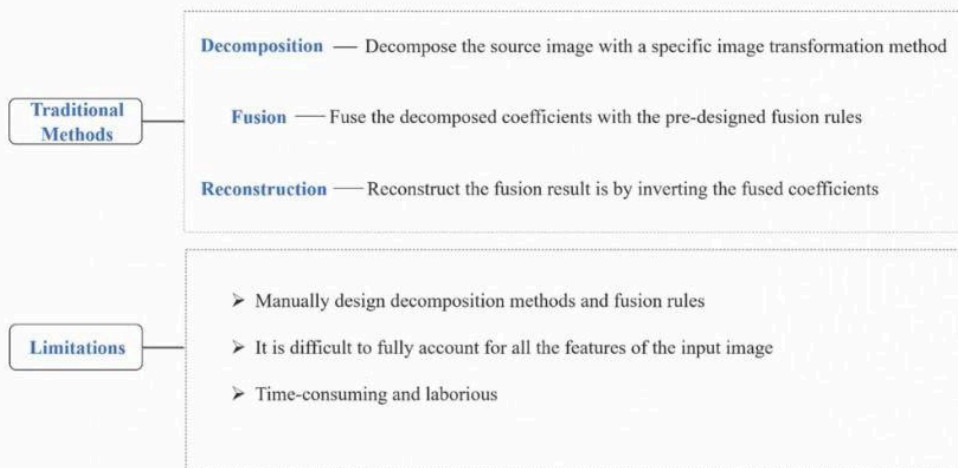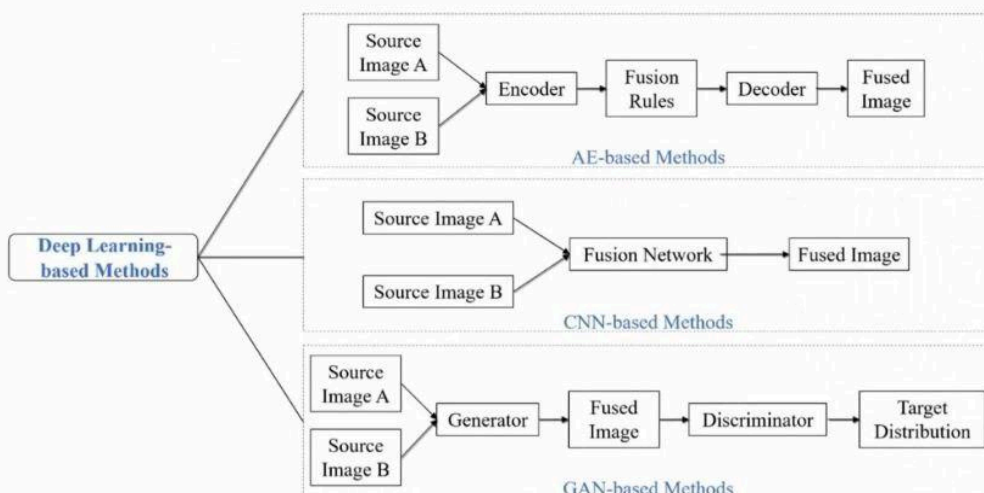


Under-exposure Image     Over-Exposure Image     Fused Image

# Infrared and Visible Image Fusion



Infrared Image     Visible Image     Fused Image

# Multi-focus Image Fusion



Left-focus Image     Right-focus Image     Fused Image

## Literature Review

**Traditional Methods**

Source Image A → Transformation/Feature Extraction → Feature Fusion → Inverse Transformation/Image reconstruction → Fused Image

Source Image B → Transformation/Feature Extraction → Feature Fusion

**Deep Learning-based Methods**

Source Image A → Fusion Network → Fused Image

Source Image B → Fusion Network

---

## Literature Review

**Traditional Methods**

**Decomposition** — Decompose the source image with a specific image transformation method

**Fusion** — Fuse the decomposed coefficients with the pre-designed fusion rules

**Reconstruction** — Reconstruct the fusion result is by inverting the fused coefficients

**Limitations**

➤ Manually design decomposition methods and fusion rules

➤ It is difficult to fully account for all the features of the input image

➤ Time-consuming and laborious

---

## Literature Review

**Deep Learning-based Methods**

Source Image A, Source Image B → Encoder → Fusion Rules → Decoder → Fused Image
*AE-based Methods*

Source Image A, Source Image B → Fusion Network → Fused Image
*CNN-based Methods*

Source Image A, Source Image B → Generator → Fused Image → Discriminator → Target Distribution
*GAN-based Methods*

Infrared Image

Significant Contrast Information

Visible Image

Rich Texture Details

Fusion Model

Fused Image

Preserving useful information from both the infrared and visible modalities.

Infrared Image

Visible Image

Fusion Network

Fused Image

Different modalities

Unified fusion network

Neglect of unique features in each source image.

Incomplete extraction of complementary information
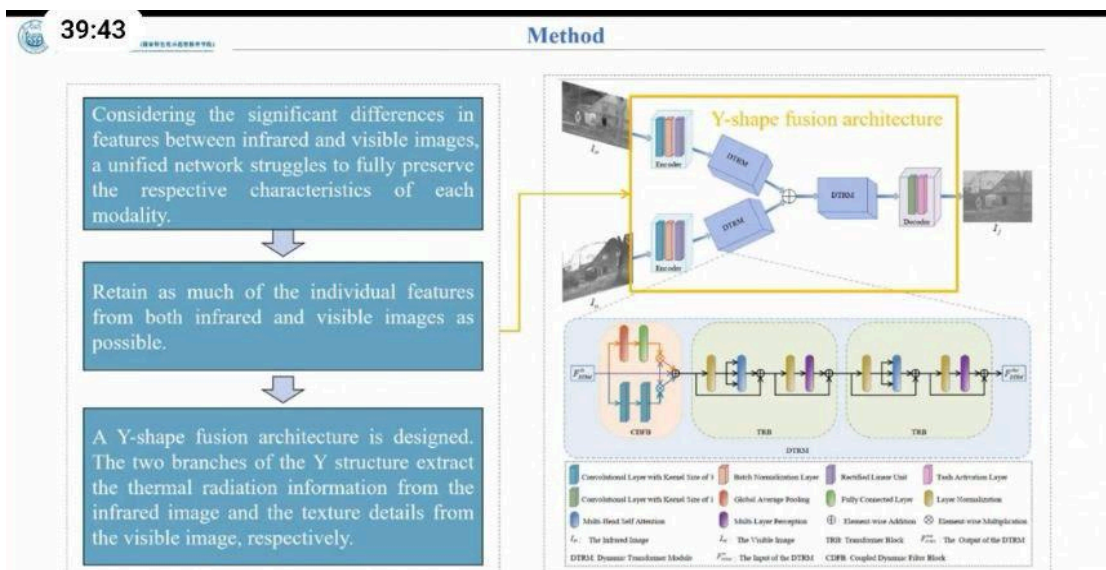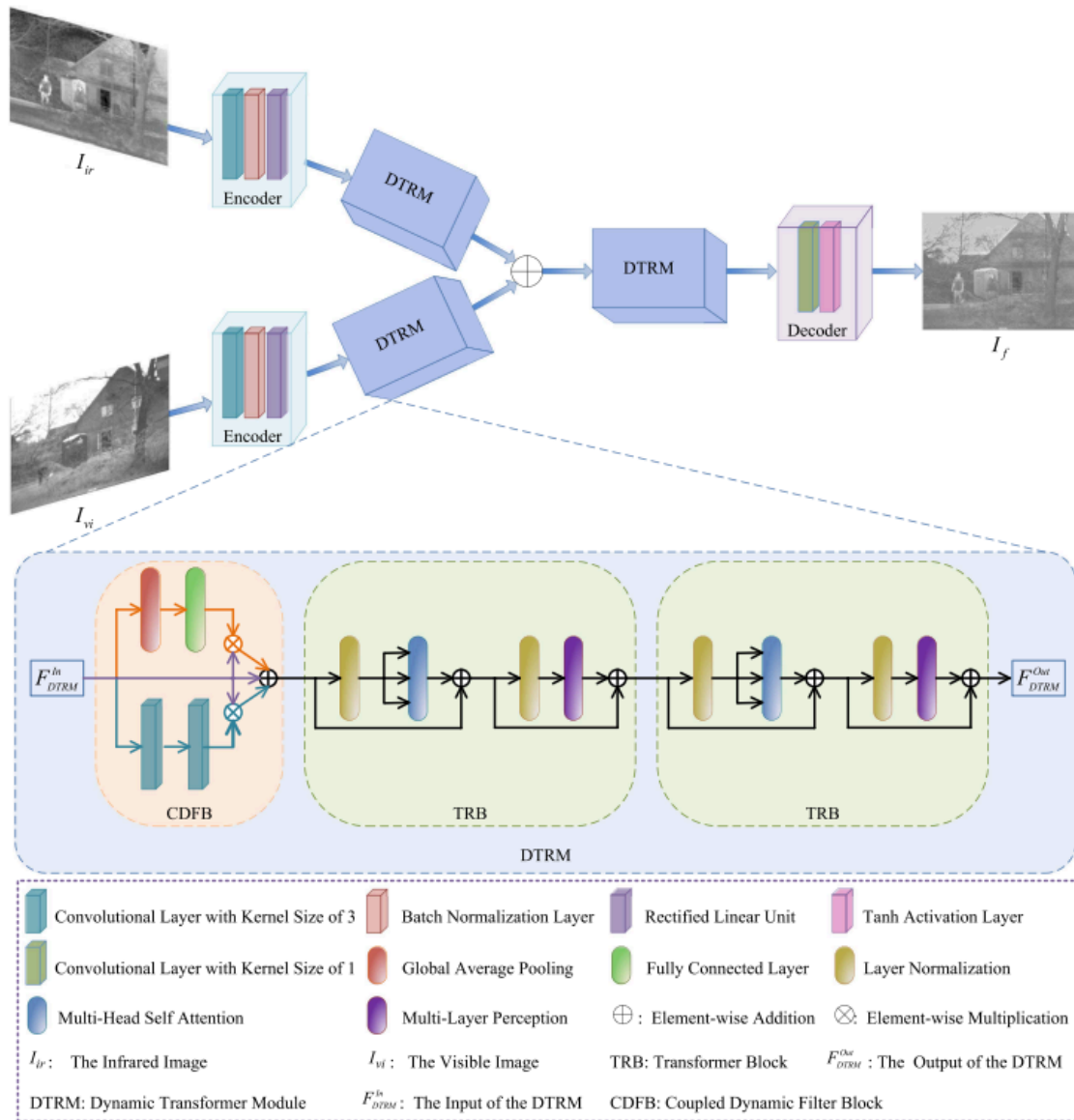
# Y-shape Dynamic Transformer (YDTR)

- **Goal**: Develop a novel infrared and visible image fusion method
- **Objective**: Generate a composite image that combines:
  - Salient target in the infrared image.
  - Texture details from the visible image.

- **Problem with existing methods**: Current deep learning based methods rely on convolutional operations, which limit global feature preservation.

- **Proposed solution**:
  - YDTR uses a dynamic Transformer module (DTRM) to capture both local features and significant context information
- **YDTR architecture** consists of:
  - **Two Y-shaped branches**: One branch extracts thermal information from the IR image, and the other extracts texture details from the visible image.
  - Each branch uses an **encoder** to capture shallow features and a **dynamic Transformer module (DTRM)** to model long-range relationships.
  - The **main path** combines these features through a DTRM and a decoder to reduce dimensions and integrate the information.

- **Loss function**: Combines two terms to enhance fusion quality:
  - Structural similarity (SSIM)
  - Spatial frequency (SF)

- **Extension**: YDTR can be extended to handle:
  - Infrared and RGB-visible images.
  - Multi-focus images.

- **Generalization**: The method demonstrates strong generalization capability without requiring fine-tuning.

IR images - Thermal radiation information

RGB images - texture details



The figure shows the overall network architecture. The top pipeline processes the infrared image $I_{ir}$ and visible image $I_{vi}$ through Encoders and DTRM modules, combined via element-wise addition, passed through another DTRM and a Decoder to produce the fused image $I_f$. The lower portion details the DTRM structure consisting of a CDFB block followed by two TRB blocks, operating on input $F_{DTRM}^{In}$ to produce output $F_{DTRM}^{Out}$.

Legend:
- Convolutional Layer with Kernel Size of 3
- Batch Normalization Layer
- Rectified Linear Unit
- Tanh Activation Layer
- Convolutional Layer with Kernel Size of 1
- Global Average Pooling
- Fully Connected Layer
- Layer Normalization
- Multi-Head Self Attention
- Multi-Layer Perception
- $\oplus$ : Element-wise Addition
- $\otimes$ : Element-wise Multiplication
- $I_{ir}$ : The Infrared Image
- $I_{vi}$ : The Visible Image
- TRB: Transformer Block
- $F_{DTRM}^{Out}$ : The Output of the DTRM
- DTRM: Dynamic Transformer Module
- $F_{DTRM}^{In}$ : The Input of the DTRM
- CDFB: Coupled Dynamic Filter Block

**Method**

If only structural loss is employed to guide network training, it's challenging to retain rich scene details.

↓

Preserve as many important details from the source images as possible.

↓

A spatial-frequency loss is designed.

$$L = L_{SSIM}(I_f, I_s) + L_{SF}(I_f, I_s)$$
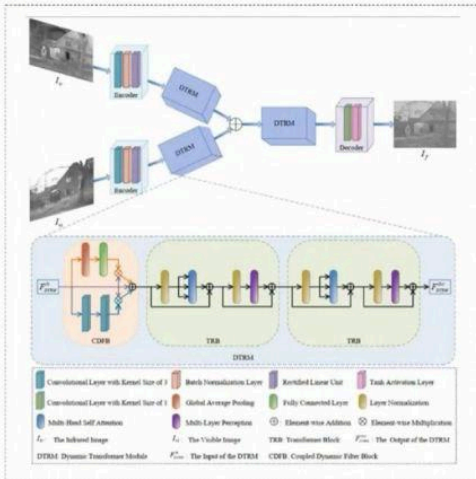
$$L_{SSIM}(I_f, I_s) = 1 - SSIM(I_f, I_s)$$

$$SSIM(I_f, I_s) = \frac{(2\mu_s\mu_f + C_1)(2\sigma_{sf} + C_2)}{(\mu_s^2 + \mu_f^2 + C_1)(\sigma_s^2 + \sigma_f^2 + C_2)}$$

$$L_{SF}(I_f, I_s) = \|SF(I_f) - SF(I_s)\|_2$$

$$SF = 1 - \sqrt{Hor^2 + Ver^2},$$

$$Hor = \sqrt{\frac{1}{HW}\sum_{i=1}^{H}\sum_{j=2}^{W}|I(i,j) - I(i,j-1)|^2}$$

$$Ver = \sqrt{\frac{1}{HW}\sum_{i=1}^{H}\sum_{j=2}^{W}|I(i,j) - I(i-1,j)|^2}.$$



**Conclusion**

The design of the **Y-shape fusion architecture** allows for the comprehensive extraction of complementary features from multi-modal images.

The design of **the dynamic transformer** facilitates the extraction of global complementary information from the input images.

Loss Functions Incorporating Structural Similarity and Spatial Frequency.

Experimental results have validated the effectiveness of the method and its practical application value.

COMPARISON WITH STATE-OF-THE-ART IMAGE FUSION ALGORITHMS

| Methods | End-to-End | Convolutional Operation | Transformer | Y-shape | SSIM Loss | SF Loss | Unsupervised | Generalization Ability |
|---|---|---|---|---|---|---|---|---|
| CNN [25] | × | ✓ | × | × | × | × | × | × |
| AUIF [26] | × | ✓ | × | × | ✓ | × | ✓ | × |
| DenseFuse [18] | × | ✓ | × | × | ✓ | × | ✓ | × |
| FusionGAN [19] | ✓ | ✓ | × | × | × | × | ✓ | × |
| AttentionFGAN [2] | ✓ | ✓ | × | ✓ | × | × | ✓ | × |
| GANMcC [24] | ✓ | ✓ | × | × | × | × | ✓ | ✓ |
| MgAN-Fuse [27] | ✓ | ✓ | × | × | × | × | ✓ | × |
| U2Fusion [20] | ✓ | ✓ | × | × | ✓ | × | ✓ | × |
| RFN-Nest [23] | ✓ | ✓ | × | × | ✓ | × | × | ✓ |
| MFE-EAG [28] | ✓ | ✓ | × | × | ✓ | × | × | ✓ |
| CSF [29] | × | ✓ | × | × | ✓ | × | ✓ | × |
| IFT [30] | × | ✓ | ✓ | × | ✓ | × | ✓ | × |
| DNDT [31] | ✓ | ✓ | ✓ | × | ✓ | × | ✓ | × |
| PPT Fusion [32] | × | × | ✓ | × | × | × | × | × |
| TGFuse [33] | ✓ | ✓ | ✓ | × | ✓ | × | ✓ | × |
| YDTR | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

# DL-Based Fusion

**Methods:**

- **CNN-based Methods**: Liu et al. [25] introduced CNNs for infrared and visible image fusion, using a Siamese network for activity level measurement. Following this, methods like **DenseFuse** [18] and **FusionGAN** [19] leveraged dense connections and GANs to enhance fusion quality.

- **Attention Mechanisms**: FusionGAN was enhanced with a **multi-scale attention mechanism** [2] to focus on discriminative regions.

- **Multi-classification and Attention**: Methods like **GANMcC** [24] and **MgAN-Fuse** [27] incorporated **multi-classification constraints** and **multi-grained attention modules** for improved fusion.

**Challenges:**

- **Limited Long-Range Context**: Convolutional operations capture local features but fail to model global context, leading to a loss of significant global features.

- **Feature Extraction Approach**: Many methods use single or parallel networks without tailoring to the specific characteristics of infrared and visible images.

**Transformer in Image Fusion:**

- The **Transformer** architecture [37], introduced for NLP, addresses CNN's limited receptive field. The **Vision Transformer (ViT)** [41] applies this to image classification, improving global feature handling.

- **Transformer-based Fusion Methods**:

    - **Multi-Scale Fusion** [30] uses a two-stage training approach.
    - **DenseNet-Transformer** [31] combines DenseNet for encoding and dual-Transformers for fusion.
    - **Patch Pyramid Transformer (PPT)** [32] uses a patch Transformer for sequence transformation and a Pyramid Transformer for feature extraction.

**Proposed Approach:**

- The paper proposes a **hybrid CNN-Transformer fusion method** to preserve both **local and global features**. This approach overcomes the limitations of existing DL-based methods, enhancing the overall fusion quality for infrared and visible images.