

# Medical Image Registration and Its Application in Retinal Images: A Review

Qiushi Nie<sup>1</sup>, Xiaoqing Zhang<sup>1,2</sup>, Yan Hu<sup>1</sup>, Mingdao Gong<sup>1</sup>,  
Jiang Liu<sup>1,3,4,5\*</sup>

<sup>1</sup>Research Institute of Trustworthy Autonomous Systems and  
Department of Computer Science and Engineering, Southern University  
of Science and Technology, Shenzhen 518055, China.

<sup>2</sup>Center for High Performance Computing and Shenzhen Key  
Laboratory of Intelligent Bioinformatics, Shenzhen Institute of Advanced  
Technology, Chinese Academy of Sciences, Shenzhen 518055, China.

<sup>3</sup>Research Institute of Trustworthy Autonomous Systems and  
Department of Computer Science and Engineering, Southern University  
of Science and Technology, Shenzhen 518055, China.

<sup>4</sup>Singapore Eye Research Institute, 169856, Singapore.

<sup>5</sup>State Key Laboratory of Ophthalmology, Optometry and Visual Science,  
Eye Hospital, Wenzhou Medical University, Wenzhou 325027, China.

\*Corresponding author(s). E-mail(s): [liuj@sustech.edu.cn](mailto:liuj@sustech.edu.cn);

Contributing authors: [12232413@mail.sustech.edu.cn](mailto:12232413@mail.sustech.edu.cn);

[xq.zhang2@siat.ac.cn](mailto:xq.zhang2@siat.ac.cn); [huy3@sustech.edu.cn](mailto:huy3@sustech.edu.cn);

[12011204@mail.sustech.edu.cn](mailto:12011204@mail.sustech.edu.cn);

## Abstract

Medical image registration is vital for disease diagnosis and treatment with its ability to merge diverse information of images, which may be captured under different times, angles, or modalities. Although several surveys have reviewed the development of medical image registration, these surveys have not systematically summarized methodologies of existing medical image registration methods. To this end, we provide a comprehensive review of these methods from traditional and deep learning-based directions, aiming to help audiences understand the development of medical image registration quickly. In particular, we review recent advances in retinal image registration at the end of each section, which has not

attracted much attention. Additionally, we also discuss the current challenges of retinal image registration and provide insights and prospects for future research.

**Keywords:** Image Registration, Medical Image, Deep Learning, Machine Learning, Retina

## 1 Introduction

Medical image registration is a fundamental step in computer-aided diagnosis (CAD) and image-guided surgical treatment, attracting much attention. It aligns multiple medical images by finding appropriate spatial transformation relationships for fusing their corresponding information, helping doctors make a more comprehensive and precise diagnosis conclusion. Particularly, these medical images may be acquired at different times, angles, and even modalities of a certain tissue or organ of the human body as input. Therefore, the nature of medical image registration is to eliminate the interference of these factors and find consistent objects or shapes for matching.

To deal with different transformation tasks in medical image registration, massive methods have been developed. They can be grouped into two types: coarse-grained global linear registration and fine-grained local elastic registration. Coarse-grained global linear registration extracts salient features of the input image pair, thereby matching these features and overcoming angular changes. Fine-grained local elastic registration performs pixel-level analysis of the input image pair after linear alignment and performs local corrections to overcome spontaneous tissue movements and deformations.

Another way to classify registration methods is according to what is used to match the images. A first and direct approach is intensity-based methods [1]. These methods consider registration as an optimization problem by iteratively disturbing the transformation parameters to maximize the pixel-wise similarity. Another early but still popular approach is feature-based methods [2], which extract manually designed features and descriptors, match them, and establish the transformation based on the matching. In contrast to intensity-based methods, feature-based methods provide more robust registration by matching salient features rather than simply comparing the pixels.

In the past decade, deep features have taken the place of handcraft features with their ability to provide learnable and, therefore, more flexible and problem-specific feature representations for registration tasks. Later, after the deep feature extractors, end-to-end registration neural networks integrated the whole registration process into a single network, applying deep-learning techniques such as convolutional neural networks (CNN), generative adversarial networks (GAN), and Transformers. Once trained, these methods can obtain registration results directly from input image pairs, thereby speeding up registration, and they have also been proven to have better registration performances.

Several reviews have been conducted on deep learning for medical image registration [3–5]. However, these papers only investigated the popular CNN-based methods at

the time, which did not mention the latest Transformer-based methods. Additionally, these works only investigated methods based on deep learning, but ignored traditional methods from the early years, which can also provide significant guidance.

Among medical images, retinal images focus on a unique part of the human body that allows for non-invasive observation of blood vessels in vivo. Retinal image registration, which combines complementary structural and functional information from the same or different modalities, is a crucial step in the process. However, in the past few years, few surveys have systematically reviewed retinal image registration. Saha *et al.* [6] and Pan *et al.* [7] reviewed retinal image registration but focused solely on registering one specific retinal modality. Moreover, they did not compare these methods with mainstream medical image registration methods.

The purpose of this paper is to review and summarize the existing medical image registration works from traditional methods and deep learning-based methods, aiming to help audiences grasp the development of medical image registration clearly. Moreover, we also survey and synthesize retinal image registration works as a particular characteristic of this review. Finally, we also highlight the current challenges of retinal image registration and discuss future research directions.

The overall organization is presented in Fig. 1: Section 2 defines the basic concepts of image registration and briefly introduces the popular retinal image modalities. Section 3 and Section 4 review the general methodology of medical image registration, the corresponding applications in retinal image registration, and comparisons between them, categorized by traditional and deep learning, respectively. Furthermore, Section 5 discusses the advantages and disadvantages of the reviewed methods, points out the current challenges, and provides potential future research directions. Finally, in Section 6, we summarize the paper.

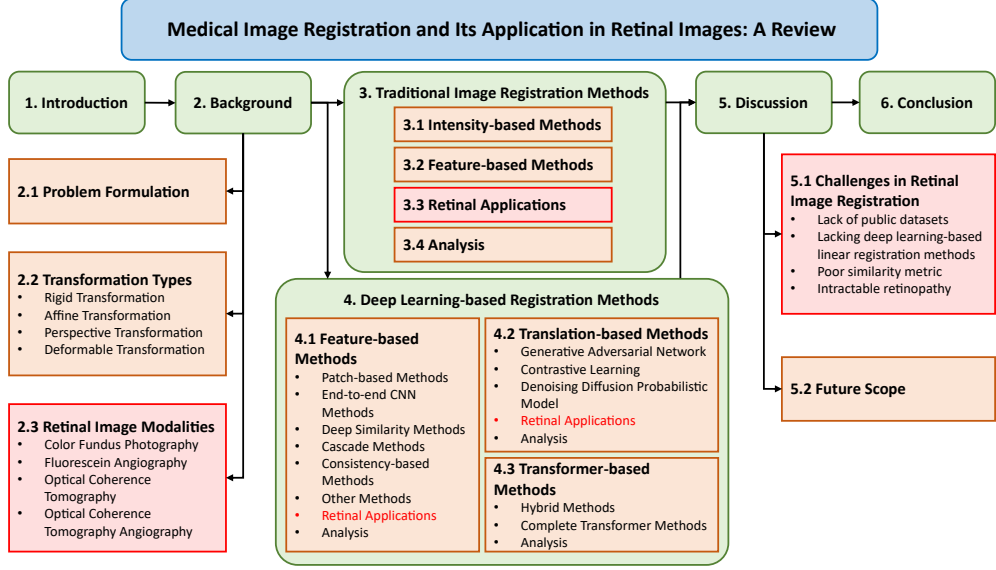
## 2 Background

### 2.1 Problem Formulation

Image registration is a fundamental task in image processing. It involves finding correspondences between two images, namely a moving and a fixed image, and establishing a transformation between them. The fixed image is used as a reference, and the goal is to transform the moving image to match the fixed image. Registration algorithms are designed to find the best transformation, denoted by  $T^*$ , that maximizes the similarity between the two images [8]. This can be achieved by maximizing the image similarity function  $\text{sim}(I_f, T(I_m))$ , where  $I_m$  and  $I_f$  are the moving and fixed images, respectively, and  $T(I_m)$  is the transformed moving image using the transform  $T$ .

### 2.2 Transformation Types

In this section, we introduce different transformation models. There are four main types of transformation: rigid, affine, perspective, and deformable. The first three are linear transformations capable of executing macro adjustments, while the fourth is non-linear and corrects local discrepancies. Fig. 2 offers a visual representation of the



**Fig. 1** Structure of our review

effects of these transformations. Although the descriptions and figures below pertain to the 2D registration, inferences can be drawn for 3D registration.

### 2.2.1 Rigid Transformation

*Rigid transformation* is a widely used simple and fast method in image processing. It consists of translation and rotation and never changes the size or shape of the original image. The model of rigid body transformation can be represented as:

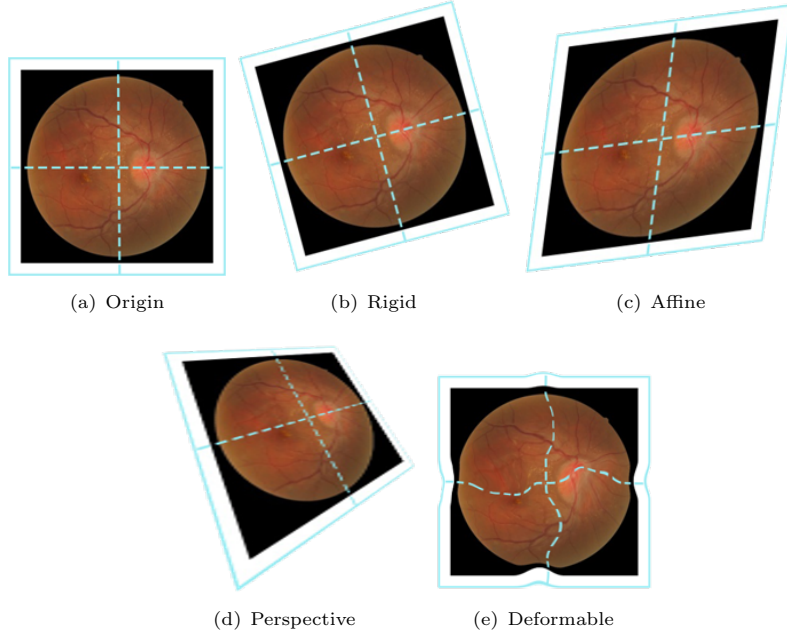
$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \mathbf{R} \begin{bmatrix} x \\ y \end{bmatrix} + \mathbf{t} \quad (1)$$

Here,  $(x, y)$  is the coordinate of the pixel in the image to be transformed, and  $(x', y')$  is the target pixel in the transformed image,  $\mathbf{R} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$  is the rotation matrix, and  $\mathbf{t} = [t_x, t_y]^T$  is the translation vector.

### 2.2.2 Affine Transformation

*Affine transformation* is realized by combining a series of atomic transformations. Based on rigid transformation, affine transformation adds scaling and shearing. The model of affine transformation can be represented as:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \mathbf{t} \quad (2)$$



**Fig. 2** The effect of different transformations

Compared to rigid transformation, affine transformation adds more freedom to the rotation matrix to handle more significant differences between images than rigid transformation. It can map straight lines to straight lines and retains the property of preserving parallelism. It maintains the parallel relationship between lines but cannot maintain the vertical relationship between lines.

### 2.2.3 Perspective Transformation

*Perspective transformation*, or projective transformation, is a more advanced form of transformation that corrects perspective distortions between images. Perspective transformation can correct more complex distortions such as foreshortening, skew, and non-parallelism.

Perspective transformation involves finding a transform matrix in homogeneous coordinates. The homogeneous coordinates use a tuple of 3 numbers  $(x_h, y_h, w_h)$  as point representation and can be translated from Cartesian coordinates  $(x_c, y_c)$  by any  $w_h \in \mathcal{R}$ ,  $x_h = w_h x_c$ ,  $y_h = w_h y_c$ . In homogeneous coordinates, the transformation can be defined as:

$$\begin{bmatrix} x' \\ y' \\ w' \end{bmatrix} = \begin{bmatrix} A & B & C \\ D & E & F \\ a & b & c \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix} \quad (3)$$

Here,  $(x, y, w)$  is the image's homogeneous coordinate to be transformed,  $(x', y', w')$  is the target coordinate in the transformed image. By setting  $w = 1$  and transform the

target  $w' = 1$ , we have the target point  $(x', y')$  back in Cartesian coordinate

$$\begin{aligned} x' &= \frac{Ax + By + C}{ax + by + c} \\ y' &= \frac{Dx + Ey + F}{ax + by + c} \end{aligned} \quad (4)$$

The perspective transformation has a straightness-preserving property, which means that a straight line is mapped to a straight line, but neither parallelism nor perpendicularity can be guaranteed.

#### 2.2.4 Deformable Transformation

*Deformable transformation* can map the shape of one image onto another through an elastic deformation model, allowing for nonlinear deformation in local regions to better adapt to different shape variations compared to rigid or affine transformation methods. In deformable transformation, we first define a deformation field that describes the amount of deformation at each pixel position and then solve for this deformation field to map one image onto another.

The model of deformable transformation can be simply represented as:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix} + \phi[x, y] \quad (5)$$

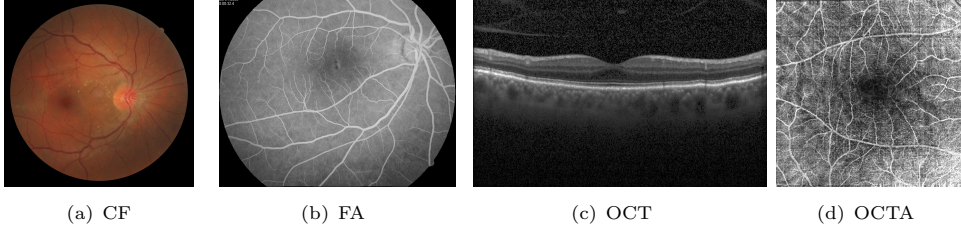
Here,  $\phi$  represents the deformation field, and  $\phi[x, y]$  represents the transformation vector  $(\Delta x, \Delta y)$  at  $(x, y)$ .

### 2.3 Retinal Image Modalities

To illustrate retinal image registration later, we introduce four commonly used techniques for photographing the eye: Color Fundus Photography (CF), Fluorescein Angiography (FA), Optical Coherence Tomography (OCT), and Optical Coherence Tomography Angiography (OCTA). These techniques provide various medical imaging tools to analyze retinal situations.

#### 2.3.1 Color Fundus Photography

Color Fundus Photography (CF) involves using a fundus camera to capture color images of the retina using white light. Equipped with a low-power microscope, the camera magnifies the view of the eye's interior surface. This technique is cost-effective and straightforward for trained professionals [9]. CF images (shown in Fig. 3(a)) contain a broader range of fundus and rich color information, making it helpful in checking the atrophy of the retina and macular. Additionally, it helps diagnose retinopathy such as diabetes retinopathy, age-related macular degeneration, and glaucoma, as well as revealing signs of systemic diseases such as diabetes and cardiovascular diseases [10].



**Fig. 3** Fundus photography examples using different imaging techniques. (a) CF from FIRE dataset [18]. (b) FA from CF-FA dataset [19]. (c) OCT from [20]. (d) OCTA from OCTA-500 dataset [21].

### 2.3.2 Fluorescein Angiography

Fluorescein Angiography (FA), shown in Fig. 3(b), involves a special dye called fluorescein and a camera to trace blood flow in the retina and choroid. It used a special dye, i.e., fluorescein, and a camera to examine blood flow in the retina and choroid. The radio-opaque dye is injected in a vein of the tester’s arm while the retina vessels are photographed by tracing the dye before and after the injection. FA can detect capillary leakage [11], aneurysm, and neovascularization. However, some people may experience discomfort after the procedure [12].

### 2.3.3 Optical Coherence Tomography

Optical Coherence Tomography (OCT) is an imaging technology that uses interference between an investigated object and a local reference signal to create high-resolution cross-sectional images and 3D scans of the retina and anterior segment [13]. Fig. 3(c) shows a cross-sectional scan of OCT. It is a non-invasive technique that enables visualization of each layer of the retina, measurement of its thickness, and provides treatment guidance for conditions such as glaucoma, Diabetic Retinopathy (DR), and Age-related Macular Degeneration (AMD). Intra-operative OCT (iOCT) is necessary for many retinal therapies, including glaucoma surgery [14] and epiretinal device implantation [15], as it provides real-time visualization of retinal layers.

### 2.3.4 Optical Coherence Tomography Angiography

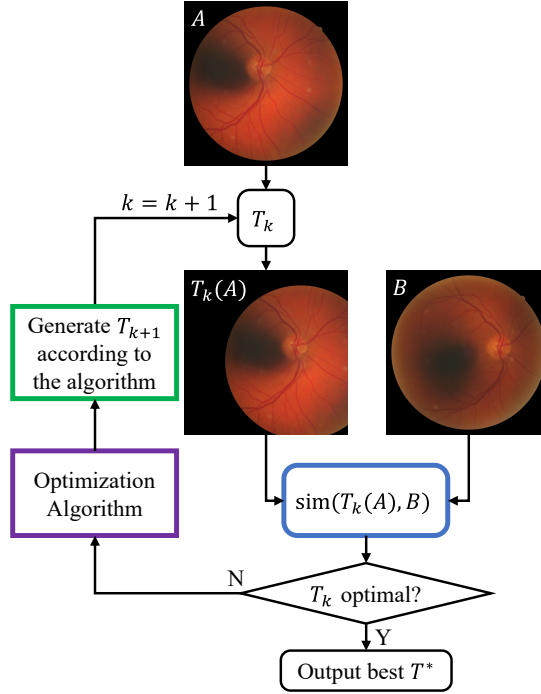
Fig. 3(d) showcases Optical Coherence Tomography Angiography (OCTA), an emerging imaging technology that builds upon OCT. OCTA captures images of the vascular network with higher resolution and a smaller view than FA, without invasiveness. Using the decorrelation signal produced by moving red blood cells, OCTA generates an image of the microvascular network. Recent studies have demonstrated the ability of OCTA to surpass the limitations of assessing blood flow in the optic nerve and help explain the vascular pathogenesis of glaucoma [16] and show impressive success in preclinical DR diagnosis [17].

### 3 Traditional Image Registration Methods

Researchers developed increasingly sophisticated algorithms and resilient features during the initial image registration phases to attain precise registration. The paper employs the phrase "traditional methods" to differentiate between techniques utilized before the advent of deep learning and those implemented after that.

#### 3.1 Intensity-based Methods

Intensity-based methods treat the problem as an iterative optimization problem. The basic steps of intensity-based registration are shown in Fig.4. Initially, a random transformation  $T_0$  is selected, and an objective function is defined to measure the similarity between the transformed image  $T_k(A)$  and the other image  $B$ . The goal is to find the optimal transformation  $T^*$  to maximize the similarity. At each step, the optimization algorithm applies a perturbation to the parameters in  $T$  based on the current similarity measure  $\text{sim}(T_k(A), T(b))$ . The process will be terminated when the similarity meets the requirement or converges with no more increase.



**Fig. 4** A general procedure of registration using iterative optimization

Researchers mainly concentrate on developing various similarity functions, including (normalized) cross-correlation (CC), (normalized) mutual information (MI), and sum of squared differences (SSD). These functions are typically calculated using the

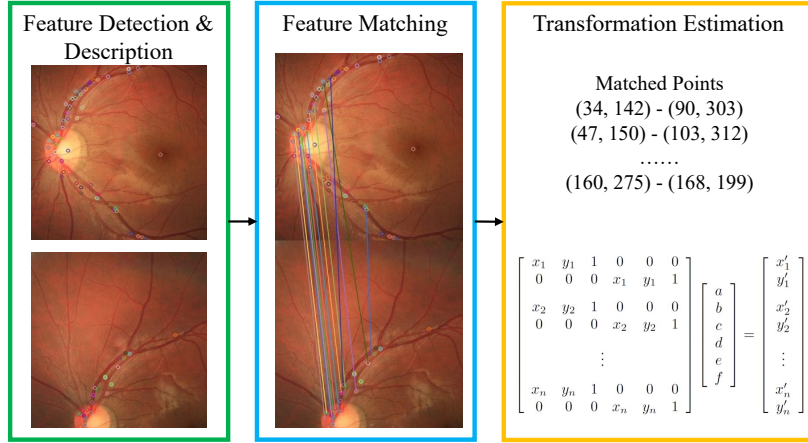


difference of each corresponding pixel of the input image pair. Among them, MI is considered the most important and widely used function. The LDDMM (Large Deformation Diffeomorphic Metric Mapping) [22] model is based on manifold learning theory and uses the Euler-Lagrange equation for optimization. It regards the image as a point on the manifold and achieves image registration by calculating the deformation between manifolds. This model can handle large deformations and maintain the nonlinear structure of the image.

Recently, there have been a few researches on intensity-based methods. Annkristin *et al.* [23] proposed a normalized gradient fields (NGF) distance measure to deal with 2D-3D image registration. To overcome the drawback that using standard similarity measures may lead to optimization problems with many local optima, Öfverstedt *et al.* [24] adopted a symmetric, intensity interpolation-free similarity measure combining intensity and spatial information. Castillo [25] proposed an intensity-based deformable image registration optimization formulation, making it easier to optimize. The similarity function is designed as a simple quadratic function formulation to be solved by straightforward coordinate descent iteration.

### 3.2 Feature-based Methods

Feature-based methods are a popular way to match images based on their correspondence. These methods focus on local structures and salient features of images rather than global information. The process is divided into three steps. First, features such as points, edges, and regions are extracted from the input images. Next, a descriptor is calculated for each feature. In the matching stage, the closest features from the two images are matched to establish potential correspondences. The idea is that the corresponding points should have very similar descriptors. Finally, the transformation parameters are estimated based on the matching results. The primary challenge is determining the most effective method for extracting and describing features. Fig. 5 illustrates the key point-based registration process.



**Fig. 5** A general procedure of keypoint-based registration

One pioneering work in feature point-based registration is the Scale-Invariant Feature Transform (SIFT) [26]. SIFT transforms the image data into scale-invariant coordinates, identifies stable key points, assigns orientations to the key points, and generates feature descriptors for each key point. The extracted feature offers invariance under variations in scale, brightness, and angles. However, this process is often computationally expensive. To address this issue, various efforts [27–31] have been made to enhance the performance and efficiency of SIFT. For instance, the Speeded Up Robust Features (SURF) [27] simplifies the filter function to reduce the dimension of descriptors and improve computational efficiency. Another method, the Oriented FAST and Rotated BRIEF (ORB) [30], integrates the FAST [32] keypoint detector and BRIEF [33] descriptor to solve the high computational cost of SIFT features and the lack of rotation invariance, scale invariance and sensitivity to noise of the BRIEF feature. As a result, ORB is capable of delivering a speedup of up to two significant figures than SIFT. Other works focus on edge and contour features, using classic edge detection [34, 35] and image segmentation [36] algorithms for feature extraction.

### 3.3 Retinal Applications

In retinal image registration, intensity-based methods are first explored. The intensity similarity metrics mentioned above, such as MI [37–39] and CC [40], are used.

The feature-based methods are more effective for retinal image registration than intensity-based methods. A popular approach is to utilize typical landmarks in retinal images. In 2003, Stewart *et al.* [41] introduced the Dual-Bootstrap Iterative Closest Point (Dual-Bootstrap ICP) algorithm for retinal image registration. This algorithm starts by matching individual vascular landmarks and aligning images based on detected blood vessel centerlines. Other studies have also utilized vascular features [42–45] and the optic disc [46] for registration purposes.

One potential solution is to enhance the capabilities of key point detectors and feature descriptors to improve performance. Ramli *et al.* [47] designed a D-Saddle detector capable of detecting feature points even in low-quality regions. Yang *et al.* [48] built upon previous work [41] to create the generalized dual-bootstrap iterative closest point (GDB-ICP), which uses better initialization, robust estimation, and strict decision criteria to align retinal images from different modalities. Chen *et al.* [49] implemented a Harris detector to identify corner points, extract partial intensity invariant feature descriptors (PIIFD), and perform bilateral matching between image pairs. Outliers are then removed, followed by applying the final transformation. Ramli *et al.* [47] improved the Saddle detector to detect feature points for low-quality regions. Gharabaghi *et al.* [50] utilized affine moment invariants (AMI) as a shape descriptor. Combining domain knowledge, SIFT and its variants are used in [51, 52]. Li *et al.* [53] introduced Orientation-independent Feature Matching (OIFM) that uses a new circular neighborhood-based feature descriptor.

### 3.4 Analysis

In the traditional registration stage, there are many applications for retinal image registration, and many registration methods directly use various retinal modalities

as evaluation indicators. Some work also specifically introduces domain knowledge based on some general methods in retinal image registration. The intensity-based approach can be sensitive to the intensity distribution when the image pair has varying illumination due to different cameras, modalities, or retinopathy-induced background changes. Feature-based methods also suffer from this problem because the feature needs descriptors. Another drawback is that most traditional registration methods take much longer for inference.

## 4 Deep Learning-based Registration Methods

Deep learning-based image segmentation has proved to be a robust tool in image segmentation since 2019 [54]. These methods can improve accuracy and efficiency by automatically learning high-level features from input images. The registration task, similar to the segmentation, has thus been developed utilizing deep learning methods. They differ from feature-based approaches by utilizing deep neural networks to replace the feature extractor, feature matching process, and transformation process. Rather than directly optimizing transformation parameters, these methods indirectly optimize registration model parameters, revealing the true essence of their effectiveness.

### 4.1 Feature-based Methods

The convolutional network (CNN) is a pioneering work in computer vision. It uses learnable convolution kernels and inductive bias, such as locality and translation equivariance, to detect learned patterns in local regions and extract high-level features. This characteristic makes CNNs especially suitable for object detection and image registration tasks, where spatial features are essential. Table 1 displays prominent works in CNN-based registration methods, which have become the most popular approach in the field since 2016.

#### 4.1.1 Patch-based Methods

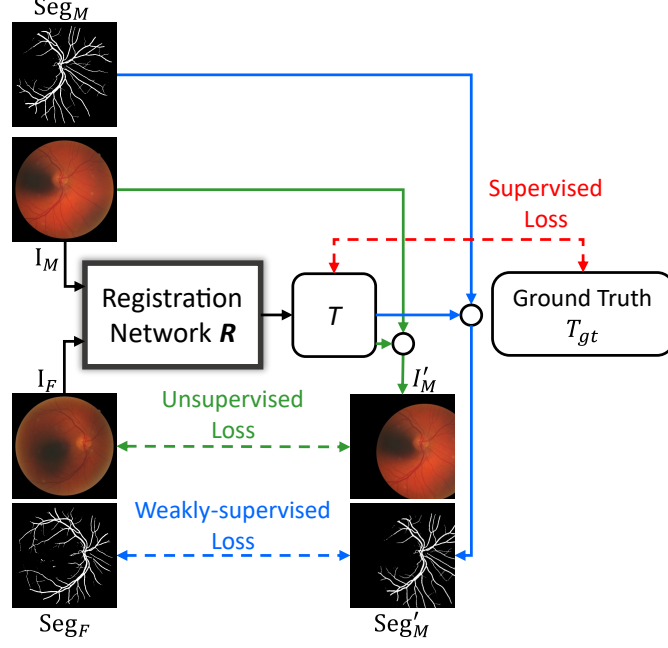
Instead of direct regression of registration parameters from the image pair, the patch-based approach was used, dividing the image into smaller patches. The patch is utilized in different ways depending on the predicted transformation type. In the case of linear transformations, the network establishes a match that can be used to derive the registration parameters. Conversely, a local displacement field is outputted and combined for nonlinear transformations. Various CNN models were proposed by Zagoruyko *et al.* [102] that output the similarity of two image patches as feature descriptors. Cao *et al.* [57] proposed a similarity-steered CNN regression architecture that estimates the displacement vectors at each corresponding location between linearly aligned brain MR image pairs. An interpolation is subsequently utilized to obtain the dense deformation. Lv *et al.* [63] divided the signal into three bins and used CNN to estimate the displacement field for abdominal motion correction throughout the respiratory cycle. However, these methods typically require an additional step of patch selection and final registration, which can be time-consuming. Additionally, generating or manually labeling ground truth can be a limiting factor.

**Table 1** Overview of feature-based image registration methods. For the supervision column, S is for supervised, W is for weakly supervised, and U is for unsupervised. For the last column, MM is for multi-modal.

Method	Year	Scene	Dimension	Modality	Type	Supervision	MM
Miao <i>et al.</i> [55]	2016	Virtual	2D/3D	X-ray/CT	Rigid	S	N
Quicksilver [56]	2017	Brain	3D	MR	Deformable	S	N
Cao <i>et al.</i> [57]	2017	Brain	3D	MR	Deformable	S	N
DIRNet [58]	2017	Digits/Heart	2D	Handwritten/MR	Deformable	U	N
Li <i>et al.</i> [59]	2017	Brain	3D	MR	Deformable	U	N
Zheng <i>et al.</i> [60]	2018	Bone	2D/3D	X-ray/CT	Rigid	S	Y
Sloan <i>et al.</i> [61]	2018	Brain	3D	MR	Rigid	S	N
AIRNet [62]	2018	Brain	3D	MR	Affine	S	N
Lv <i>et al.</i> [63]	2018	Abdominal	3D	MR	Deformable	S	N
Hu <i>et al.</i> [64]	2018	Prostate gland	3D	MR/US	Deformable	W	Y
Jiang <i>et al.</i> [65]	2018	Chest	3D	CT	Deformable	U	N
Li <i>et al.</i> [66]	2018	Brain	3D	MR	Deformable	U	N
BIRNet [67]	2019	Brain	3D	MR	Deformable	S	N
DeepAtlas [68]	2019	Knee/Brain	3D	MR	Deformable	W	N
DLIR [69]	2019	Heart/Chest	3D	MR/CT	Affine/Deformable	U	N
VTN [70]	2019	Liver/Brain	3D	CT/MR	Affine/Deformable	U	N
Zhao <i>et al.</i> [71]	2019	Liver/Brain	3D	CT/MR	Affine/Deformable	U	N
VoxelMorph [72]	2019	Brain	3D	MR	Deformable	W/U	N
VoxelMorph-diff [73]	2019	Brain	3D	MR	Deformable	U	N
Dual-PRNet [74]	2019	Brain	3D	MR	Deformable	U	N
DeepFLASH [75]	2020	Synthetic/Brain	2D/3D	Eye/MR	Deformable	S	N
Mansilla <i>et al.</i> [76]	2020	Chest	2D	X-ray	Deformable	W	N
Mok <i>et al.</i> [77]	2020	Brain	3D	MR	Deformable	U	N
CycleMorph [78]	2021	Face/Brain	2D/3D	Expression/MR	Deformable	U	N
DeepSim [79]	2021	Brain/Cell	3D	MR/EM	Deformable	W/U	N
Mok <i>et al.</i> [80]	2022	Brain	3D	MR	Deformable	U	N
Dual-PRNet <sup>++</sup> [81]	2022	Brain	3D	MR	Deformable	U	N
Tran <i>et al.</i> [82]	2022	Liver/Brain	3D	CT/MRI	Deformable	U	N
IMSE [83]	2023	Brain	2D/3D	CT/MR	Deformable	U	Y
Zhang <i>et al.</i> [84]	2023	Liver	3D	US	Deformable	U	N
AMNet [85]	2023	Brain	3D	MR	Deformable	U	N
DeepSPa [86]	2019	Retina	2D	CF / FA / OCT	Affine	U	Y
Zhang <i>et al.</i> [87]	2019			CF / FA	Deformable	W	Y
Silva <i>et al.</i> [88]	2020			CF / FA / IR	Affine	S	Y
Wang <i>et al.</i> [89]	2020			CF / IR	Perspective	S	Y
Tian <i>et al.</i> [90]	2020			CF / OCT	Deformable	U	Y
Zou <i>et al.</i> [91]	2020			CF	Deformable	U	N
Wang <i>et al.</i> [92]	2021			CF / FA / IR	Perspective	S	Y
Zhang <i>et al.</i> [93]	2021			CF / FA / IR	Affine/Deformable	S	Y
Sui <i>et al.</i> [94]	2021			MSI	Deformable	W	N
An <i>et al.</i> [95]	2022			CF / FA / IR	Rigid	U	Y
Benvenuto <i>et al.</i> [96]	2022			CF	Deformable	U	N
Lopez <i>et al.</i> [97]	2022			OCTA	Deformable	U	N
Rivas <i>et al.</i> [98]	2022			CF	Similarity	S	N
Kim <i>et al.</i> [99]	2022			CF	Perspective	S	N
Rivas <i>et al.</i> [100]	2023			OCT	Affine + Z-axis	U	N

#### 4.1.2 End-to-end CNN Methods

Due to increased computing power, supervised end-to-end networks are developed for direct registration. The ground truth is obtained by traditional algorithms or manual labels. A general end-to-end deep learning registration method framework is shown in Figure 6. Miao *et al.* [55] employed 2D/3D CNN regressors to estimate the rigid transformation parameters in real time directly. Quicksilver [56] divides the 3D brain MRI into 3D patches due to the limitation of GPU memory, but it can directly predict the deformation field for the input patches. To improve the performance of supervised methods, Chee *et al.* [62] leveraged unlabelled data to generate a synthetic dataset and further trained the AIRNet (affine image registration network) based on it. BIRNet



**Fig. 6** The overall framework for end-to-end deep learning-based medical image registration methods. Red, blue, and green lines denote the supervised, weakly-supervised, and unsupervised training strategies, respectively. The small circles denote performing spatial transformation with the predicted transform  $T$  using STN [101].

[67] was proposed as a hierarchical dual-supervised fully convolutional neural network based on U-Net [103] in the following year, with a loss function designed as a combination of the difference in image intensity and the difference of predicted displacement and ground truth displacement in each layer of U-Net’s decoder. Wang *et al.* [75] introduced a low dimensional Fourier representation of diffeomorphic transformations to improve training and inference efficiency.

Weakly supervised registration methods take advantage of additional semantic information to ensure meaningful registration, and they also overcome the challenge of the unavailability of ground truth transformation. These methods utilize extra information, such as anatomical segmentation, to perform registration. Hu *et al.* [64] proposed a weakly supervised registration network for multi-modal 3D prostate gland images using the ground truth segmentation label of the gland and other anatomical landmarks. Xu *et al.* [68] proposed a deep learning framework named DeepAtlas that jointly learns networks for image registration and segmentation, which are trained alternately, complementing each other to achieve better results with only a few labels for segmentation.

Unsupervised methods have also been researched to eliminate any ground truth labels further. Spatial transformer layer (STL) [101], a differentiable module that can warp the input image, has been the foundation of many unsupervised registration methods. STL enables obtaining the transformed moving image in a differentiable

manner, which allows applying conventional similarity measurement between the transformed and fixed images during training as the loss function. In 2017, DIRNet [58] was introduced as the first end-to-end unsupervised deformable registration network by adopting STL. Later, VoxelMorph [72] was proposed as a U-Net-based network that achieved faster run time and better performance than traditional iterative-based methods, with only unsupervised training. Auxiliary anatomical segmentation can be optionally used in a weakly-supervised setting. In their following work, Dalca *et al.* [73] further adopted a probabilistic generative model to provide diffeomorphic guarantees. Dual-PRNet [74] extended VoxelMorph [72] by incorporating a pyramid registration module that uses multi-level context information and sequentially warps the convolutional features. Dual-PRNet<sup>++</sup> [81] further enhanced the PR module in Dual-PRNet by computing correlation features and using residual convolutions.

#### 4.1.3 Deep Similarity Methods

In deep learning, pixel-based similarity metrics like MSE and NCC are commonly employed. However, these metrics may sometimes encounter difficulties when dealing with low-intensity contrast or noise. To address these issues, deep similarity methods have been developed, which utilize a custom similarity measure. For example, DeepSim [79] utilizes semantic information extracted by a pre-trained feature extractor in a segmentation network to construct a semantic similarity metric. This specialized metric allows the network to learn and adapt to dataset-specific features, improving low-quality image performance. IMSE [83] takes it a step further with a self-supervised approach to train a modality-independent evaluator using a new data augmentation technique called shuffle remap, which can provide style enhancement. The evaluator then serves as a multi-modal similarity estimator to train a multi-modal registration network.

#### 4.1.4 Cascade Methods

Cascade methods are inspired by traditional iterative registration. The cascade architecture, namely, stacking networks in series, can provide progressive registration in a coarse-to-fine manner. DLIR [69] implemented a cascade architecture by stacking an affine network followed by multiple deformable networks, with each network being trained sequentially with the weights of previous networks fixed. In contrast, Zhao *et al.* [70, 71] proposed a recursive cascade architecture similar to DLIR but much more sophisticated. They jointly trained their cascade networks to learn progressive alignments more effectively.

#### 4.1.5 Consistency-based Methods

Consistency-based methods add consistency constraints based on the property of registration or transformation. In 2020, Mok *et al.* tackled the challenge of the deformable transformation’s invertibility by introducing a swift and symmetrical diffeomorphic image registration approach [77]. The network was trained with an inverse-consistency constraint, enabling it to learn bidirectional transformations to the mean shape of two input images to produce topology-preserving and inverse-consistent transformations.

The following year, Kim *et al.* proposed CycleMorph [78], which utilized cycle consistency as an additional constraint to enhance topology preservation and reduce folding issues. To register images X to Y and Y to X, the method employs two CNNs,  $G_X$  and  $G_Y$ . The warped images from both networks are used as image pairs and sent to the networks themselves to ensure they can be returned to their original state, maximizing the similarity between the original image and the reversed image.

#### 4.1.6 Other Methods

However, with the development of novel architectures, the parameters increased significantly, so it is harder to achieve real-time registration without high computing power. Tran *et al.* [82] tried to solve this problem via knowledge distillation. They transferred meaningful knowledge of distilled deformations from a pre-trained high-performance network (teacher network) to a fast, lightweight network (student network). After training, only the lightweight student network is used during the inference, allowing the model to achieve fast inference time using only a common CPU.

#### 4.1.7 Retinal Applications

Utilizing retinal landmarks inspired researchers to develop new techniques for detecting these landmarks automatically. In particular, [86] used handcrafted features, while [98, 99] employed CNNs. Lee *et al.* employed a CNN to classify patches of various step patterns with intensity changes. On the other hand, Rivas *et al.* [98] used a CNN to produce a heatmap of blood vessels and bifurcations and applied the maxima detection and feature matching method RANSAC [104] during testing. Similarly, Kim *et al.* [99] used a vessel segmentation network and a joint detection network to identify vascular landmark points for registration. The SIFT algorithm [26] was then used to compute descriptors based on the region around these points. Benvenuto *et al.* [96] used Isotropic Undecimated Wavelet Transform, which segments blood vessels and ocular shape. Based on the segmentation, the registration network adopted from U-Net is trained to perform registration. This year, Rivas *et al.* [100] also explored deep learning registration methods for OCT 3D Scan. They first performed affine alignment on a 2D projection and then z-axis registration based on layer segmentation.

Recent studies have explored the potential of end-to-end methods utilizing innovative network architectures. Silva *et al.* [88] utilized a VGG 16 feature extractor, correlation matrix, and regression network to replicate the traditional steps of feature-based registration, including feature extraction, matching, and registration transform. They evaluated the model’s effectiveness on a multi-modal retinal dataset. Meanwhile, Tian *et al.* [90] improved U-Net [103] using image pyramid as multi-scale input and introduced a new edge similarity loss calculated via the correlation between gradients of the fixed and moved image. Still based on U-Net, Sui *et al.* [94] further sent an image pyramid of the original image and the ground truth vessel map into each layer of encoder and decoder, respectively.

It is worth noting that Wang *et al.* [87, 89, 92, 93, 95] made significant strides in multi-modal retinal image registration. They began with a deformable registration model, which included a vessel segmentation network and a deformation field estimation network, in [87]. In their later work [89], they adapted the vessel segmentation



network from [87] and further used a pre-trained SuperPoint [105] model for feature detection and description, along with an outlier rejection network for perspective registration. They utilized this three-stage methodology (segmentation, detection and description, and outlier rejection) in their subsequent research. The advantage of this approach is that it bypasses the intensity gap of different modalities, but the downside is its high complexity. They improved the segmentation network with pixel-adaptive convolution in [92]. In [93], they utilized perspective registration as coarse registration and added a deformable framework for fine alignment, achieving remarkable accuracy. Finally, their latest work [95] made the three-stage methodology a self-supervised one.

#### 4.1.8 Analysis

Limited by insufficient computing power in the early stage, patch-based methods appeared earlier. As computing power and network diversity increase, it has been able to take the entire image and even 3D images as input and perform feature extraction and matching tasks integrated and simultaneously. During this period, thanks to the rapid advances in deep learning architecture, the feature extraction module for image registration has also been steadily developing along with the trend. In addition, many novel constraints have been proposed. These constraints exploit some fundamental properties of the registration or transformation, resulting in a more suitable shape for the output transformation.

Deep learning-based retinal image registration did not appear until 2019. There are two types of approaches to this method. The first type employs mainstream registration methods which utilize advanced network architectures. Although these methods achieve exceptional performance, they rely heavily on the design of network architecture rather than domain knowledge. The second type of approach focuses on solving the registration problem in a more domain-specific manner. These methods extract or leverage critical features such as vessel segmentation or vascular junctions for later registration.

We have observed that the diversity of retinal image registration works is notably lower when compared to other medical registration works. This can be attributed to a few factors, including imaging principles and targets. Regarding imaging principles, CT and MR use X-rays and magnetic fields. These imaging techniques can produce high tissue contrast while maintaining the intensity consistency of different images. However, when it comes to retinal image registration, CF and FA images are commonly used, which rely solely on white light illumination for imaging. This imaging principle, coupled with the natural movement of the subject’s eyeballs, can lead to differences in brightness when taking multiple shots. The second factor is the imaging target. CT and MR are typically used for imaging the chest, abdomen, and brain, with many features such as multiple organs, tissues, and brain regions. In contrast, retinal images mainly focus on the blood vessel tree and optic disc, with features that have relatively low discrimination. As a result, learning robust features in retinal image registration automatically proves to be a more challenging task.

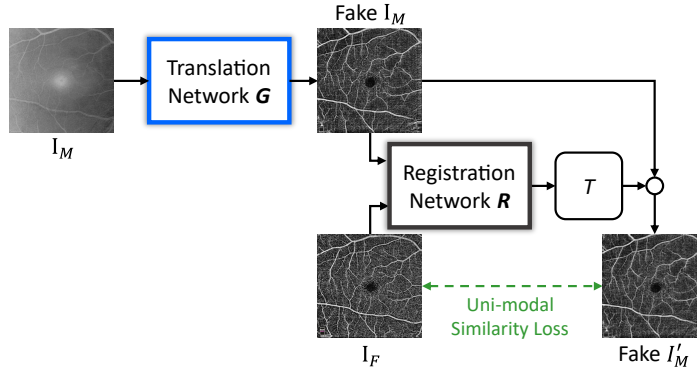


**Table 2** Overview of Translation-based image registration methods

Methodology	Method	Year	Scene	Dimension	Modality	Type
GAN	Mahapatra <i>et al.</i> [106]	2018	Retina/Heart	2D	CF/FA/MR	Deformable
	Qin <i>et al.</i> [107]	2019	Lung/Brain	2D	CT/MR	Deformable
	Xu <i>et al.</i> [108]	2020	Kidney/Abdomen	3D	CT/MR	Deformable
	Han <i>et al.</i> [109]	2022	Brain	3D	MR/CT	Deformable
	MedRegNet [110]	2022	Retina	2D	CF/FAF/FAG	Perspective
Contrastive Learning	Casamitjana <i>et al.</i> [111]	2021	Brain	3D	Histology/MRI	Deformable
	Chen <i>et al.</i> [112]	2022	Thorax/Abdomen/Lung	3D	CT/MRI	Deformable
DDPM	DiffuseMorph [113]	2022	Face/Brain	2D/3D	Expression/MR	Deformable

## 4.2 Translation-based Methods

Multi-modal image registration can be complex as it involves aligning images of varying modalities with unique intensity distributions. This can pose a challenge for uni-modal methods. However, an innovative solution to this issue is to leverage image translation techniques. This solution transforms the multi-modal registration problem into a more straightforward uni-modal registration problem, as depicted in Fig. 7. Table 2 highlights the top translation-based registration algorithms available.

**Fig. 7** The overall framework for translation-based methods.

### 4.2.1 Generative Adversarial Network

Generative adversarial network (GAN) [114] consists of two sub-networks, a generator and a discriminator, trained in a game-theoretic setting to generate synthetic data indistinguishable from the actual data. The generator generates synthetic samples while the discriminator attempts to differentiate between natural and synthetic samples, and the training process continues until the generated samples are indistinguishable from actual ones.

Mahapatra *et al.* [106] used a GAN to generate the registered image with the identical distribution of the moving image and the deformation field. They also ensured that the structure of the generated image matched that of the reference image through structural similarity loss. Qin *et al.* [107] proposed a method of decomposing images

into a latent shape space and a separate latent appearance space for both modalities, which was used to learn a bi-directional registration function.

CycleGAN [115], which is based on GAN, enables image-to-image (i2i) translation using unpaired images. It employs a cycle consistency loss to ensure the reconstructed images are consistent with the original input images. Several multi-modal registration methods [108, 109] used CycleGAN as the primary network for image translation. Xu *et al.* [108] introduced two additional losses to enforce structural similarity between the translated and authentic images. They also jointly trained translated uni-modal and multi-modal streams to complement each other. Han *et al.* [109] implemented image synthesis in both directions and predicted the associated uncertainty, providing information used in the fusion of the two direction estimations.

#### 4.2.2 Contrastive Learning

Contrastive learning defines positive and negative samples, and the goal is to learn a representation space where positive samples are close to each other while negative ones are far away. A recent study by Park *et al.* [116] explored integrating contrastive learning into image translation by introducing an extra loss called PatchNCE to naive GAN. This loss encourages generated output patches to be closer to their corresponding image patches than random ones. Casamitjana *et al.* [111] used the PatchNCE loss to train an i2i translation network for transferring source images to the desired target domain. They then applied an independently trained intra-modality registration network in the target domain to predict the deformation field. Building on this work, Chen *et al.* [112] proposed an end-to-end architecture that jointly trained the registration and translation network without requiring a discriminator.

#### 4.2.3 Denoising Diffusion Probabilistic Model

A new generative model called the denoising diffusion probabilistic model (DDPM) [117] was recently introduced. This model is designed to learn the Markov transformation from a simple Gaussian distribution to the actual data distribution. DDPM has been shown to generate images of higher quality than GAN [118]. Additionally, Kim *et al.* [113] developed DiffuseMorph, the first and currently only registration network based on diffusion. The network estimates the score function by adding a diffusion network before a standard registration network, and it even shows the image registration trajectory by scaling the conditional score. However, unlike translation between modalities, DiffuseMorph constructs the score function directly between the input image pair.

#### 4.2.4 Retinal Applications

Although there have been many studies on image-to-image translation in various retinal modalities [129, 130], few studies on retinal image registration use translation-based techniques. While MedRegNet [110] is the only method that utilizes CycleGAN [115] as an image translation tool, it is solely used as a multi-modal retinal data generator. On the other hand, [87, 89, 92, 93, 95] mentioned in the previous part can also be considered as translation-based methods when dealing with multi-modal data.

**Table 3** Overview of Transformer-based image registration methods

Method	Year	Scene	Dimension	Modality	Type	MM
ViT-V-Net [119]	2021	Brain	3D	MR	Deformable	N
Zhang <i>et al.</i> [120]	2021	Brain	3D	MR	Deformable	N
C2FViT [121]	2022	Brain	3D	MR	Affine	N
TransMorph [122]	2022	Brain/Heart	3D	MRI/XCAT/CT	Affine/Deformable	N
TD-Net [123]	2022	Brain	3D	MR	Deformable	N
Wang <i>et al.</i> [124]	2022	Brain	3D	MR	Deformable	N
XMorpher [125]	2022	Heart	3D	CT	Deformable	N
Swin-VoxelMorph [126]	2022	Brain	3D	MR	Deformable	N
TransMatch [127]	2023	Brain	3D	MR	Deformable	N
ModeT [128]	2023	Brain	3D	MR	Deformable	N

They employ image segmentation to obtain blood vessel segmentation maps, which can be viewed as converting different modalities into a single "mask" modality.

#### 4.2.5 Analysis

Translation-based methods have been shown to effectively solve multi-modal problems by transforming image pairs into the same modality, thus reducing registration difficulty. Recently, registration methods using different types of generative networks have emerged with the updated generative network models. While GAN can produce good results in modality translation, its training process is quite complex, requiring more time-consuming manual hyperparameter adjustment for the generator and discriminator. In the past, contrastive learning accounted for half of the unsupervised field, but it typically requires large amounts of high-quality data for training. Diffusion is another recently emerging image generation method that can produce quite realistic effects, but its applicability in registration still needs to be explored. Unfortunately, research progress in this area is limited, and the need for more public datasets for retinal image registration may be one of the reasons.

### 4.3 Transformer-based Methods

Recently, Google explored a way to use pure transformer architecture in vision tasks, known as Vision Transformer (ViT) [131], outperforming existing CNN methods' performance. ViTs split the image into patches and treat them like tokens as in an NLP application, which has led to its successful application in various computer vision tasks, including image registration. Table 3 overviews transformer-based image registration methods.

#### 4.3.1 Hybrid Methods

In the beginning, researchers attempted to integrate Transformers into CNN-based models. Chen *et al.* [119] pioneered the use of ViT on high-level features extracted from convolutional layers of moving and fixed images. Building on this approach, Song [123] proposed TD-Net, which utilizes multiple Transformer blocks for down-sampling. Conversely, Zhang *et al.* [120] introduced a dual Transformer network comprised of two branches - intra-image and inter-image - with Transformers embedded in both

branches to enhance features, similar to the approach in [119]. Wang *et al.* [124] enhanced the UNet [103] architecture for registration by introducing a bi-level connection and a unique Transformer block. TransMorph [122] was then proposed as a hybrid Transformer-ConvNet model, utilizing Swin Transformers [132] in the encoder and convolutional layers in the decoder. The authors demonstrated that positional embedding could be disregarded, leading to a flatter loss landscape for registration. The following year, Chen *et al.* [127] proposed TransMatch, emphasizing the importance of inter-image feature matching. They employed a Transformer-based encoder and matched regions using their new local window cross-attention module. Recently, Wang *et al.* [128] introduced a motion decomposition transformer (ModeT) based on a multi-head neighborhood attention mechanism which can model multiple motion modalities.

### 4.3.2 Complete Transformer Methods

An alternative method involves integrating a complete Transformer architecture into the network. In a recent study by Shi *et al.* [125], a unique X-shaped transformer architecture called XMorpher was introduced. The researchers incorporated cross-attention between two feature extraction branches and a window size constraint to enhance information exchange and locality of the network. Another study, Swin-VoxelMorph [126], utilized a fully Swin Transformer-based 3D Swin-UNet and a bidirectional constraint to optimize both forward and inverse transformation. To fill the gap of affine image registration, Mok *et al.* [121] proposed a Coarse-to-Fine Vision Transformer (C2FViT), a pure transformer architecture. The researchers transformed image pairs into small-to-large resolutions and passed them through different stages of ViT to achieve the desired results.

### 4.3.3 Analysis

In CNN, convolution is usually performed within a local region, whereas the Transformer’s self-attention mechanism allows for global information exchange. Therefore, integrating the Transformer into the registration network primarily aims to leverage its ability to construct global information, further enhancing the feature extraction. Additionally, there is a growing trend towards designing pure Transformer architectures, demonstrating outstanding performance in visual tasks. However, the attention mechanism should be used in a more registration-specific manner. To address this problem, cross-attention Transformers are being explored at both the feature extraction and feature matching stages.

It is worth noting that nearly all Transformer-based registration methods are applied in the Brain MRI dataset, so there is no related work on retinal registration. This phenomenon may be caused by access to a more significant number of MR public datasets and an enormous amount of data. At the same time, ViT requires a significantly more significant amount of data to achieve better results than the general CNN network.

## 5 Discussion

### 5.1 Challenges in Retinal Image Registration

#### 5.1.1 Lack of public datasets

In artificial intelligence, many tasks rely on competition and public evaluation challenges to make progress. These challenges offer a comprehensive and impartial platform for researchers to compare the performance, computation time, and robustness of newly designed algorithms. The Learn2Reg challenge, for example, recently focused on registering medical image modalities commonly used in brain, abdomen, and thorax imaging [133]. The currently available datasets for retinal image registration are recorded in Table 4. It has not formed enough public retinal image datasets for each modality, nor has there been any competition.

#### 5.1.2 Lacking deep learning-based linear registration methods

Based on the articles using deep learning we reviewed, we calculated the proportion of different transformation types used in general registration methodology and the proportion of each type specifically in the retinal image registration task, as shown in Fig. 8. We found that over 80% of the works on general methodology employ non-linear transformation. However, for retinal applications, linear transformation is the most common method. This is because retinal images are primarily captured from a particular area of the retina, and when registering different images, the contents of these different visual field areas are generally registered. On the other hand, most of the organs captured by commonly used modalities are 3D images and have no difference in field of view. The main differences are in position, inevitable elastic deformation of the organs, or differences caused by different collection objects. Therefore, the registration methods mainly focus on correcting non-linear deformations. Most will directly use traditional algorithms to perform rigid registration for position differences. Because these mainstream registration methods mainly focus on non-linear registration, linear registration tasks for retinal images are not popular.

#### 5.1.3 Poor similarity metric

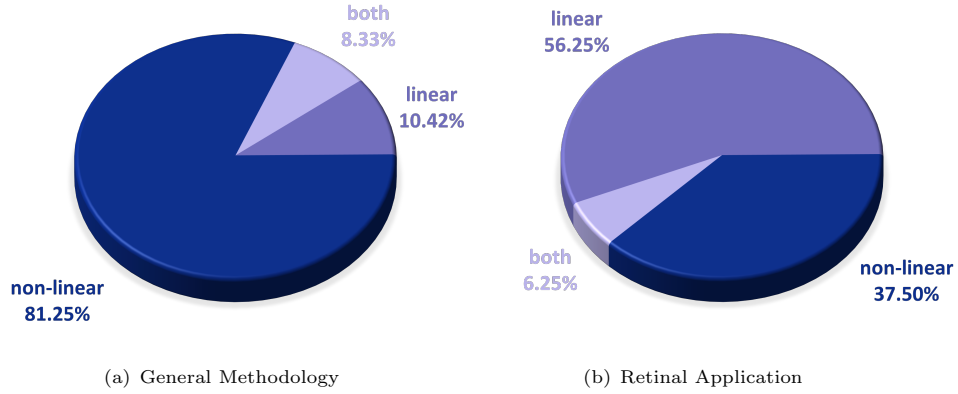
Similar metrics are always used to optimize the registration network in an unsupervised manner or to evaluate the quality of registration. The key technical challenge in medical image registration is selecting and designing the most effective similarity measurement. In uni-modal image registration, brightness changes may present the most significant difficulty. One of the main obstacles in multi-modal image registration is that images from different modalities hold different resolutions, contrast, and luminosity. A newly designed similarity metric or a completely different technical route for multi-modal image registration is urgently needed.

#### 5.1.4 Intractable retinopathy

In clinical treatments, most patients suffer from eye retinopathy, so their retina may be severely damaged. Small bulges, swellings, or blood may cover the normal fundus and

**Table 4** Public retinal image registration datasets

Dataset	Source	Camera Specs.	Format	Modality	Resolution	Size (pairs)	Ground truth
FIRE [18]	Papageorgiou Hospital, Aristotle University of Thessaloniki, Greece	Nidek AFC-210 fundus camera	JPG	CF	2912×2912	134	Control Points
FLoRI21 [134]	RECOVERY study [135]	Optos California and 200Tx cameras	TIFF	UWF FA	3900×3072	15	Control Points
CF-FFA [19]	Unknown	Unknown	JPG	CF & FFA	720×576	60	None

**Fig. 8** Deep learning-based method using different transformation types

have a bad effect on photography. Meanwhile, some diseases may change the structure of the retina. Most samples in the public datasets are retinal images of ordinary people. However, when it is used in clinical diagnosis, some patients' retinas are likely to have retinopathy. In this situation, a network trained using normal images cannot work well.

## 5.2 Future Scope

In this era of large models, we may anticipate a general large model for registration to emerge soon. With the ability to use human-marked point pairs or corresponding mask areas as registration prompts, this model will be trained on higher quality, broader types, and more extensive image registration datasets, allowing for better generalization.

There are still plenty of areas to explore regarding retinal image registration. With multiple image modalities, there is a pressing need for multi-modal image registration. To address this, translation-based methods and disentangling representation methods may be another new exploration. Interestingly, we have not seen any pioneers attempting Transformer-based retinal registration methods, which could lead to even greater accuracy.

Moreover, data scarcity remains a significant challenge, but we may overcome it through data generation or transfer learning. For instance, we can supplement the dataset with image pairs obtained through random translation, rotation, and brightness and contrast enhancement using retinal images from other datasets. When employing transfer learning, we can train on endoscopic images from other parts of the human body or manually generate virtual datasets and fine-tune them for retinal image registration.

## 6 Conclusion

This paper thoroughly analyzed medical image registration, focusing on its applications in retinal imaging. Our review compares general medical image registration techniques and their adaptation to retinal imaging, highlights gaps in current research, and gives a little advice on avenues for future exploration. We also evaluated state-of-the-art medical image registration methods and weighed the advantages and disadvantages of each. Lastly, we identified challenges specific to retinal registration and discussed potential opportunities for further advancement.

## Declarations

**Abbreviations.** CAD: Computer-Aided Diagnosis; CNN: Convolutional Neural Network; GAN: Generative Adversarial Network; CF: Color Fundus photography; FA: Fluorescein Angiography; OCT: Optical Coherence Tomography; OCTA: Optical Coherence Tomography Angiography; DR: Diabetic Retinopathy; AMD: Age-related Macular Degeneration; CC: Cross Correlation; MI: Mutual Information; SSD: Sum of Squared Difference; DDPM: Denoising Diffusion Probabilistic Model; ViT: Vision Transformer.

**Availability of data and materials.** All relevant data and material are presented in the main paper.

**Competing interests.** The authors declare that they have no competing interests.

**Funding.** This work was supported in part by General Program of National Natural Science Foundation of China (82102189 and 82272086), and Guangdong Provincial Department of Education (SJZLGC202202).

**Authors' contributions.** All the authors make substantial contribution in this manuscript. QN, YH, and MG participated in writing the first draft of the paper. Then the paper was revised carefully and completed the final manuscript by QN and XZ. JL supervised the study. All the authors have read and approved the final manuscript.

**Acknowledgement.** Not applicable.

## References

- [1] Oliveira, F.P., Tavares, J.M.R.: Medical image registration: a review. *Computer methods in biomechanics and biomedical engineering* **17**(2), 73–93 (2014)

- [2] Zitova, B., Flusser, J.: Image registration methods: a survey. *Image and vision computing* **21**(11), 977–1000 (2003)
- [3] Boveiri, H.R., Khayami, R., Javidan, R., Mehdizadeh, A.: Medical image registration using deep neural networks: a comprehensive review. *Computers & Electrical Engineering* **87**, 106767 (2020)
- [4] Haskins, G., Kruger, U., Yan, P.: Deep learning in medical image registration: a survey. *Machine Vision and Applications* **31**, 1–18 (2020)
- [5] Bharati, S., Mondal, M., Podder, P., Prasath, V.: Deep learning for medical image registration: A comprehensive review. *arXiv preprint arXiv:2204.11341* (2022)
- [6] Saha, S.K., Xiao, D., Bhuiyan, A., Wong, T.Y., Kanagasingam, Y.: Color fundus image registration techniques and applications for automated analysis of diabetic retinopathy progression: A review. *Biomedical Signal Processing and Control* **47**, 288–302 (2019)
- [7] Pan, L., Chen, X.: Retinal oct image registration: methods and applications. *IEEE Reviews in Biomedical Engineering* (2021)
- [8] Khalifa, F., Beache, G.M., Gimel’farb, G., Suri, J.S., El-Baz, A.S.: State-of-the-art medical image registration methodologies: A survey. *Multi Modality State-of-the-Art Medical Image Segmentation and Registration Methodologies: Volume 1*, 235–280 (2011)
- [9] Besenczi, R., Tóth, J., Hajdu, A.: A review on automatic analysis techniques for color fundus photographs. *Computational and Structural Biotechnology Journal* **14**, 371–384 (2016) <https://doi.org/10.1016/j.csbj.2016.10.001>
- [10] Abràmoff, M.D., Garvin, M.K., Sonka, M.: Retinal imaging and image analysis. *IEEE Reviews in Biomedical Engineering* **3**, 169–208 (2010) <https://doi.org/10.1109/RBME.2010.2084567>
- [11] Baek, J., Lee, M.Y., Kim, B., Choi, A., Kim, J., Kwon, H., Jeon, S.: Ultra-widefield fluorescein angiography findings in patients with macular edema following cataract surgery. *Ocular Immunology and Inflammation* **29**(3), 610–614 (2021) <https://doi.org/10.1080/09273948.2019.1691739> <https://doi.org/10.1080/09273948.2019.1691739>. PMID: 31850812
- [12] A.D.A.M. Medical Encyclopedia: Fluorescein angiography. Medlineplus. Retrieved Oct. 12, 2022
- [13] Podoleanu, A.G.: Optical coherence tomography. *Journal of microscopy* **247**, 209–219 (2012) <https://doi.org/10.1111/j.1365-2818.2012.03619.x>



- [14] Ang, B.C.H., Lim, S.Y., Dorairaj, S.: Intra-operative optical coherence tomography in glaucoma surgery-a systematic review. *Eye (London, England)* **34**1, 168–177 (2020) <https://doi.org/10.1038/s41433-019-0689-3>
- [15] Grewal, D.S., Carrasco-Zevallos, O., Gunther, R., Izatt, J.A., Toth, C.A., Hahn, P.: Intra-operative microscope-integrated swept-source optical coherence tomography guided placement of argus ii retinal prosthesis. *Acta ophthalmologica* **95**, 431–432 (2017) <https://doi.org/10.1111/aos.13123>
- [16] Werner, A.C., Shen, L.Q.: A review of OCT angiography in glaucoma. *Seminars in Ophthalmology* **34**(4), 279–286 (2019) <https://doi.org/10.1080/08820538.2019.1620807> <https://doi.org/10.1080/08820538.2019.1620807>. PMID: 31158045
- [17] Shaikh, N.F., Vohra, R., Balaji, A., Azad, S.V., Chawla, R., Kumar, V., Venkatesh, P., Kumar, A.: Role of optical coherence tomography-angiography in diabetes mellitus: Utility in diabetic retinopathy and a comparison with fluorescein angiography in vision threatening diabetic retinopathy **69**(11), 3218–3224 (2021) [https://doi.org/10.4103/ijo.IJO\\_1267\\_21](https://doi.org/10.4103/ijo.IJO_1267_21)
- [18] Hernandez-Matas, C., Zabulis, X., Triantafyllou, A., Anyfanti, P., Douma, S., Argyros, A.A.: FIRE: Fundus Image Registration Dataset. <https://doi.org/10.35119/maio.v1i4.42>
- [19] Alipour, S.H.M., Rabbani, H., Akhlaghi, M.R.: Diabetic retinopathy grading by digital curvelet transform. *Computational and Mathematical Methods in Medicine* **2012** (2012) <https://doi.org/10.1155/2012/761901>
- [20] Mooney, P.: Retinal OCT Images (optical coherence tomography). Retrieved Nov. 23, 2022 from <https://www.kaggle.com/datasets/paultimothymooney/kermany2018> (2017)
- [21] Li, M., Chen, Y., Xie, K., Yuan, S., Chen, Q.: OCTA-500. <https://doi.org/10.1109/TMI.2020.2992244> . <https://dx.doi.org/10.1109/TMI.2020.2992244>
- [22] Beg, M.F., Miller, M.I., Trounev, A., Younes, L.: Computing large deformation metric mappings via geodesic flows of diffeomorphisms. *International journal of computer vision* **61**, 139–157 (2005)
- [23] Lange, A., Heldmann, S.: Multilevel 2d-3d intensity-based image registration. In: Špiclin, J. Žigaand McClelland, Kybic, J., Goksel, O. (eds.) *Biomedical Image Registration*, pp. 57–66. Springer, Cham (2020)
- [24] Öfverstedt, J., Lindblad, J., Sladoje, N.: Fast and robust symmetric image registration based on distances combining intensity and spatial information. *IEEE Transactions on Image Processing* **28**(7), 3584–3597 (2019) <https://doi.org/10.1109/TIP.2019.2899947>

- [25] Castillo, E.: Quadratic penalty method for intensity-based deformable image registration and 4dct lung motion recovery. *Medical Physics* **46**(5), 2194–2203 (2019) <https://doi.org/10.1002/mp.13457> <https://aapm.onlinelibrary.wiley.com/doi/pdf/10.1002/mp.13457>
- [26] Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *International journal of computer vision* **60**, 91–110 (2004)
- [27] Bay, H., Tuytelaars, T., Van Gool, L.: Surf: Speeded up robust features. *Lecture notes in computer science* **3951**, 404–417 (2006)
- [28] Ke, Y., Sukthankar, R.: Pca-sift: A more distinctive representation for local image descriptors. In: *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, vol. 2, p. (2004). IEEE
- [29] Tola, E., Lepetit, V., Fua, P.: A fast local descriptor for dense matching. In: *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8 (2008). IEEE
- [30] Rublee, E., Rabaud, V., Konolige, K., Bradski, G.: Orb: An efficient alternative to sift or surf. In: *2011 International Conference on Computer Vision*, pp. 2564–2571 (2011). <https://doi.org/10.1109/ICCV.2011.6126544>
- [31] Cai, G.-R., Jodoin, P.-M., Li, S.-Z., Wu, Y.-D., Su, S.-Z., Huang, Z.-K.: Perspective-sift: An efficient tool for low-altitude remote sensing image registration. *Signal Processing* **93**(11), 3088–3110 (2013)
- [32] Rosten, E., Drummond, T.: Machine learning for high-speed corner detection. In: *Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7–13, 2006. Proceedings, Part I* 9, pp. 430–443 (2006). Springer
- [33] Calonder, M., Lepetit, V., Strecha, C., Fua, P.: Brief: Binary robust independent elementary features. In: *Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010. Proceedings, Part IV* 11, pp. 778–792 (2010). Springer
- [34] Canny, J.: A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence* (6), 679–698 (1986)
- [35] Marr, D., Hildreth, E.: Theory of edge detection. *Proceedings of the Royal Society of London. Series B. Biological Sciences* **207**(1167), 187–217 (1980)
- [36] Pal, N.R., Pal, S.K.: A review on image segmentation techniques. *Pattern recognition* **26**(9), 1277–1294 (1993)

- [37] Legg, P.A., Rosin, P.L., Marshall, D., Morgan, J.E.: Improving accuracy and efficiency of mutual information for multi-modal retinal image registration using adaptive probability density estimation. *Computerized Medical Imaging and Graphics* **37**(7-8), 597–606 (2013)
- [38] Reel, P.S., Dooley, L.S., Wong, K.C.P., Börner, A.: Robust retinal image registration using expectation maximisation with mutual information. In: 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 1118–1122 (2013). IEEE
- [39] Reel, P.S., Dooley, L.S., Wong, K.P., Börner, A.: Enhanced retinal image registration accuracy using expectation maximisation and variable bin-sized mutual information. In: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6632–6636 (2014). IEEE
- [40] Cideciyan, A.V.: Registration of ocular fundus images: an algorithm using cross-correlation of triple invariant image descriptors. *IEEE Engineering in Medicine and Biology Magazine* **14**(1), 52–58 (1995) <https://doi.org/10.1109/51.340749>
- [41] Stewart, C.V., Tsai, C.-L., Roysam, B.: The dual-bootstrap iterative closest point algorithm with application to retinal image registration. *IEEE Transactions on Medical Imaging* **22**(11), 1379–1394 (2003) <https://doi.org/10.1109/TMI.2003.819276>
- [42] Guo, X., Hsu, W., Lee, M.L., Wong, T.Y.: A tree matching approach for the temporal registration of retinal images. In: 2006 18th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'06), pp. 632–642 (2006). IEEE
- [43] Zheng, Y., Hunter, A.A., Wu, J., Wang, H., Gao, J., Maguire, M.G., Gee, J.C.: Landmark matching based automatic retinal image registration with linear programming and self-similarities. In: Székely, G., Hahn, H.K. (eds.) *Information Processing in Medical Imaging*, pp. 674–685. Springer, Berlin, Heidelberg (2011)
- [44] Zheng, Y., Daniel, E., Hunter, A.A., Xiao, R., Gao, J., Li, H., Maguire, M.G., Brainard, D.H., Gee, J.C.: Landmark matching based retinal image alignment by enforcing sparsity in correspondence matrix. *Medical Image Analysis* **18**(6), 903–913 (2014) <https://doi.org/10.1016/j.media.2013.09.009> . Sparse Methods for Signal Reconstruction and Medical Image Analysis
- [45] S. Hervella, Rouco, J., Novo, J., Ortega, M.: Multimodal registration of retinal images using domain-specific landmarks and vessel enhancement. *Procedia Computer Science* **126**, 97–104 (2018) <https://doi.org/10.1016/j.procs.2018.07.213> . Knowledge-Based and Intelligent Information & Engineering Systems: Proceedings of the 22nd International Conference, KES-2018, Belgrade, Serbia
- [46] Koukounis, D., Nicholson, L., Bull, D.R., Achim, A.: Retinal image registration

based on multiscale products and optic disc detection. In: 2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society, pp. 6242–6245 (2011). IEEE

- [47] Ramli, R., Idris, M.Y.I., Hasikin, K., A Karim, N.K., Abdul Wahab, A.W., Ahmady, I., Ahmady, F., Kadri, N.A., Arof, H.: Feature-based retinal image registration using d-saddle feature. *Journal of healthcare engineering* **2017** (2017)
- [48] Yang, G., Stewart, C.V., Sofka, M., Tsai, C.-L.: Registration of challenging image pairs: Initialization, estimation, and decision. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **29**(11), 1973–1989 (2007) <https://doi.org/10.1109/TPAMI.2007.1116>
- [49] Chen, J., Tian, J., Lee, N., Zheng, J., Smith, R.T., Laine, A.F.: A partial intensity invariant feature descriptor for multimodal retinal image registration. *IEEE Transactions on Biomedical Engineering* **57**(7), 1707–1718 (2010) <https://doi.org/10.1109/TBME.2010.2042169>
- [50] Gharabaghi, S., Daneshvar, S., Sedaaghi, M.H.: Retinal image registration using geometrical features. *Journal of digital imaging* **26**, 248–258 (2013)
- [51] Ghassabi, Z., Shanbehzadeh, J., Mohammadzadeh, A., Ostadzadeh, S.S.: Colour retinal fundus image registration by selecting stable extremum points in the scale-invariant feature transform detector. *IET image processing* **9**(10), 889–900 (2015)
- [52] Saha, S.K., Xiao, D., Frost, S., Kanagasalingam, Y.: A two-step approach for longitudinal registration of retinal images. *Journal of Medical Systems* **40**(277) (2016) <https://doi.org/10.1007/s10916-016-0640-0>
- [53] Li, Q., Li, S., Wu, Y., Guo, W., Qi, S., Huang, G., Chen, S., Liu, Z., Chen, X.: Orientation-independent feature matching (oifm) for multimodal retinal image registration. *Biomedical Signal Processing and Control* **60**, 101957 (2020)
- [54] Hesamian, M.H., Jia, W., He, X., Kennedy, P.: Deep learning techniques for medical image segmentation: achievements and challenges. *Journal of digital imaging* **32**, 582–596 (2019)
- [55] Miao, S., Wang, Z.J., Zheng, Y., Liao, R.: Real-time 2d/3d registration via cnn regression. In: 2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI), pp. 1430–1434 (2016). IEEE
- [56] Yang, X., Kwitt, R., Styner, M., Niethammer, M.: Quicksilver: Fast predictive image registration – a deep learning approach. *NeuroImage* **158**, 378–396 (2017) <https://doi.org/10.1016/j.neuroimage.2017.07.008>

- [57] Cao, X., Yang, J., Zhang, J., Nie, D., Kim, M., Wang, Q., Shen, D.: Deformable image registration based on similarity-steered cnn regression. In: Medical Image Computing and Computer Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I 20, pp. 300–308 (2017). Springer
- [58] De Vos, B.D., Berendsen, F.F., Viergever, M.A., Staring, M., Išgum, I.: End-to-end unsupervised deformable image registration with a convolutional neural network. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3, pp. 204–212 (2017). Springer
- [59] Li, H., Fan, Y.: Non-rigid image registration using fully convolutional networks with deep self-supervision. arXiv preprint arXiv:1709.00799 (2017)
- [60] Zheng, J., Miao, S., Jane Wang, Z., Liao, R.: Pairwise domain adaptation module for cnn-based 2-d/3-d registration. *Journal of Medical Imaging* **5**(2), 021204–021204 (2018)
- [61] Sloan, J.M., Goatman, K.A., Siebert, J.P.: Learning rigid image registration-utilizing convolutional neural networks for medical image registration (2018)
- [62] Chee, E., Wu, Z.: Airnet: Self-supervised affine registration for 3d medical images using neural networks. arXiv preprint arXiv:1810.02583 (2018)
- [63] Lv, J., Yang, M., Zhang, J., Wang, X.: Respiratory motion correction for free-breathing 3d abdominal mri using cnn-based image registration: a feasibility study. *The British journal of radiology* **91**(xxxx), 20170788 (2018)
- [64] Hu, Y., Modat, M., Gibson, E., Li, W., Ghavami, N., Bonmati, E., Wang, G., Bandula, S., Moore, C.M., Emberton, M., Ourselin, S., Noble, J.A., Barratt, D.C., Vercauteren, T.: Weakly-supervised convolutional neural networks for multimodal image registration. *Medical Image Analysis* **49**, 1–13 (2018) <https://doi.org/10.1016/j.media.2018.07.002>
- [65] Jiang, P., Shackelford, J.A.: Cnn driven sparse multi-level b-spline image registration. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 9281–9289 (2018)
- [66] Li, H., Fan, Y.: Non-rigid image registration using self-supervised fully convolutional networks without training data. In: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018), pp. 1075–1078 (2018). IEEE
- [67] Fan, J., Cao, X., Yap, P.-T., Shen, D.: Birnet: Brain image registration using

- dual-supervised fully convolutional networks. *Medical Image Analysis* **54**, 193–206 (2019) <https://doi.org/10.1016/j.media.2019.03.006>
- [68] Xu, Z., Niethammer, M.: Deepatlas: Joint semi-supervised learning of image registration and segmentation. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II* **22**, pp. 420–429 (2019). Springer
  - [69] De Vos, B.D., Berendsen, F.F., Viergever, M.A., Sokooti, H., Staring, M., Išgum, I.: A deep learning framework for unsupervised affine and deformable image registration. *Medical image analysis* **52**, 128–143 (2019)
  - [70] Zhao, S., Lau, T., Luo, J., Chang, E.I.-C., Xu, Y.: Unsupervised 3d end-to-end medical image registration with volume tweening network. *IEEE Journal of Biomedical and Health Informatics* **24**(5), 1394–1404 (2020) <https://doi.org/10.1109/JBHI.2019.2951024>
  - [71] Zhao, S., Dong, Y., Chang, E.I.-C., Xu, Y.: Recursive cascaded networks for unsupervised medical image registration. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019)
  - [72] Balakrishnan, G., Zhao, A., Sabuncu, M.R., Guttag, J., Dalca, A.V.: Voxel-morph: a learning framework for deformable medical image registration. *IEEE transactions on medical imaging* **38**(8), 1788–1800 (2019)
  - [73] Dalca, A.V., Balakrishnan, G., Guttag, J., Sabuncu, M.R.: Unsupervised learning of probabilistic diffeomorphic registration for images and surfaces. *Medical image analysis* **57**, 226–236 (2019)
  - [74] Hu, X., Kang, M., Huang, W., Scott, M.R., Wiest, R., Reyes, M.: Dual-stream pyramid registration network. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II*, pp. 382–390 (2019). Springer
  - [75] Wang, J., Zhang, M.: Deepflash: An efficient network for learning-based medical image registration. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020)
  - [76] Mansilla, L., Milone, D.H., Ferrante, E.: Learning deformable registration of medical images with anatomical constraints. *Neural Networks* **124**, 269–279 (2020)
  - [77] Mok, T.C., Chung, A.: Fast symmetric diffeomorphic image registration with convolutional neural networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4644–4653 (2020)

- [78] Kim, B., Kim, D.H., Park, S.H., Kim, J., Lee, J.-G., Ye, J.C.: Cyclemorph: cycle consistent unsupervised deformable image registration. *Medical image analysis* **71**, 102036 (2021)
- [79] Czolbe, S., Krause, O., Feragen, A.: Deepsim: Semantic similarity metrics for learned image registration. In: *Medical Imaging with Deep Learning* (2021). PMLR
- [80] Mok, T.C., Chung, A.C.: Robust image registration with absent correspondences in pre-operative and follow-up brain mri scans of diffuse glioma patients. In: *International MICCAI Brainlesion Workshop*, pp. 231–240 (2022). Springer
- [81] Kang, M., Hu, X., Huang, W., Scott, M.R., Reyes, M.: Dual-stream pyramid registration network. *Medical Image Analysis* **78**, 102379 (2022)
- [82] Tran, M.Q., Do, T., Tran, H., Tjiputra, E., Tran, Q.D., Nguyen, A.: Light-weight deformable registration using adversarial learning with distilling knowledge. *IEEE transactions on medical imaging* **41**(6), 1443–1453 (2022)
- [83] Kong, L., Qi, X.S., Shen, Q., Wang, J., Zhang, J., Hu, Y., Zhou, Q.: Indescribable multi-modal spatial evaluator. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9853–9862 (2023)
- [84] Zhang, J., Fu, T., Li, J., Xiao, D., Fan, J., Lin, Y., Song, H., Ji, F., Yang, M., Yang, J.: An alternately optimized generative adversarial network with texture and content constraints for deformable registration of 3d ultrasound images. *Physics in Medicine and Biology* (2023)
- [85] Che, T., Wang, X., Zhao, K., Zhao, Y., Zeng, D., Li, Q., Zheng, Y., Yang, N., Wang, J., Li, S.: Amnet: Adaptive multi-level network for deformable registration of 3d brain mr images. *Medical Image Analysis* **85**, 102740 (2023)
- [86] Lee, J.A., Liu, P., Cheng, J., Fu, H.: A deep step pattern representation for multimodal retinal image registration. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019)
- [87] Zhang, J., An, C., Dai, J., Amador, M., Bartsch, D.-U., Borooah, S., Freeman, W.R., Nguyen, T.Q.: Joint vessel segmentation and deformable registration on multi-modal retinal images based on style transfer. In: *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 839–843 (2019). IEEE
- [88] Silva, T., Chew, E., Hotaling, N., Cukras, C.: Deep-learning based multi-modal retinal image registration for longitudinal analysis of patients with age-related macular degeneration. *Biomedical Optics Express* **12** (2020) <https://doi.org/10.1364/BOE.408573>
- [89] Wang, Y., Zhang, J., An, C., Cavichini, M., Jhingan, M., Amador-Patarroyo,

- M.J., Long, C.P., Bartsch, D.-U.G., Freeman, W.R., Nguyen, T.Q.: A segmentation based robust deep learning framework for multimodal retinal image registration. In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1369–1373 (2020). IEEE
- [90] Tian, Y., Hu, Y., Ma, Y., Hao, H., Mou, L., Yang, J., Zhao, Y., Liu, J.: Multi-scale u-net with edge guidance for multimodal retinal image deformable registration. In: 2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC), pp. 1360–1363 (2020). IEEE
- [91] Zou, B., He, Z., Zhao, R., Zhu, C., Liao, W., Li, S.: Non-rigid retinal image registration using an unsupervised structure-driven regression network. *Neurocomputing* **404**, 14–25 (2020)
- [92] Wang, Y., Zhang, J., Cavichini, M., Bartsch, D.-U.G., Freeman, W.R., Nguyen, T.Q., An, C.: Robust content-adaptive global registration for multimodal retinal images using weakly supervised deep-learning framework. *IEEE Transactions on Image Processing* **30**, 3167–3178 (2021)
- [93] Zhang, J., Wang, Y., Dai, J., Cavichini, M., Bartsch, D.-U.G., Freeman, W.R., Nguyen, T.Q., An, C.: Two-step registration on multi-modal retinal images via deep neural networks. *IEEE Transactions on Image Processing* **31**, 823–838 (2021)
- [94] Sui, X., Zheng, Y., Jiang, Y., Jiao, W., Ding, Y.: Deep multispectral image registration network. *Computerized Medical Imaging and Graphics* **87**, 101815 (2021)
- [95] An, C., Wang, Y., Zhang, J., Nguyen, T.Q.: Self-supervised rigid registration for multimodal retinal images. *IEEE Transactions on Image Processing* **31**, 5733–5747 (2022)
- [96] Benvenuto, G.A., Colnago, M., Casaca, W.: Unsupervised deep learning network for deformable fundus image registration. In: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1281–1285 (2022). IEEE
- [97] López-Varela, E., Novo, J., Fernández-Vigo, J.I., Moreno-Morillo, F.J., Ortega, M.: Unsupervised deformable image registration in a landmark scarcity scenario: choroid octa. In: International Conference on Image Analysis and Processing, pp. 89–99 (2022). Springer
- [98] Rivas-Villar, D., Hervella, Á.S., Rouco, J., Novo, J.: Color fundus image registration using a learning-based domain-specific landmark detection methodology. *Computers in Biology and Medicine* **140**, 105101 (2022)



- [99] Kim, G.Y., Kim, J.Y., Lee, S.H., Kim, S.M.: Robust detection model of vascular landmarks for retinal image registration: A two-stage convolutional neural network. *BioMed Research International* **2022** (2022)
- [100] Rivas-Villar, D., Motschi, A.R., Pircher, M., Hitzenberger, C.K., Schranz, M., Roberts, P.K., Schmidt-Erfurth, U., Bogunović, H.: Automated inter-device 3d oct image registration using deep learning and retinal layer segmentation. *Biomedical Optics Express* **14**(7), 3726–3747 (2023)
- [101] Jaderberg, M., Simonyan, K., Zisserman, A., kavukcuoglu, k.: Spatial transformer networks. In: Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*, vol. 28. Curran Associates, Inc., ??? (2015)
- [102] Zagoruyko, S., Komodakis, N.: Learning to compare image patches via convolutional neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4353–4361 (2015)
- [103] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III* 18, pp. 234–241 (2015). Springer
- [104] Fischler, M.A., Bolles, R.C.: Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM* **24**(6), 381–395 (1981)
- [105] DeTone, D., Malisiewicz, T., Rabinovich, A.: Superpoint: Self-supervised interest point detection and description. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 224–236 (2018)
- [106] Mahapatra, D., Antony, B., Sedai, S., Garnavi, R.: Deformable medical image registration using generative adversarial networks. In: *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pp. 1449–1453 (2018). IEEE
- [107] Qin, C., Shi, B., Liao, R., Mansi, T., Rueckert, D., Kamen, A.: Unsupervised deformable registration for multi-modal images via disentangled representations. In: *Information Processing in Medical Imaging: 26th International Conference, IPMI 2019, Hong Kong, China, June 2–7, 2019, Proceedings* 26, pp. 249–261 (2019). Springer
- [108] Xu, Z., Luo, J., Yan, J., Pulya, R., Li, X., Wells, W., Jagadeesan, J.: Adversarial uni-and multi-modal stream networks for multimodal image registration. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part*

- [109] Han, R., Jones, C.K., Lee, J., Wu, P., Vagdargi, P., Uneri, A., Helm, P.A., Luciano, M., Anderson, W.S., Siewerdsen, J.H.: Deformable mr-ct image registration using an unsupervised, dual-channel network for neurosurgical guidance. *Medical image analysis* **75**, 102292 (2022)
- [110] Santarossa, M., Kilic, A., Burchard, C., Schmarje, L., Zelenka, C., Reinhold, S., Koch, R., Roider, J.: Medregnet: unsupervised multimodal retinal-image registration with gans and ranking loss. In: *Medical Imaging 2022: Image Processing*, vol. 12032, pp. 321–333 (2022). SPIE
- [111] Casamitjana, A., Mancini, M., Iglesias, J.E.: Synth-by-reg (SbR): Contrastive learning for synthesis-based registration of paired images. In: *Simulation and Synthesis in Medical Imaging*, pp. 44–54. Springer, ??? (2021). [https://doi.org/10.1007/978-3-030-87592-3\\_5](https://doi.org/10.1007/978-3-030-87592-3_5)
- [112] Chen, Z., Wei, J., Li, R.: Unsupervised Multi-Modal Medical Image Registration via Discriminator-Free Image-to-Image Translation (2022)
- [113] Kim, B., Han, I., Ye, J.C.: Diffusemorph: unsupervised deformable image registration using diffusion model. In: *European Conference on Computer Vision*, pp. 347–364 (2022). Springer
- [114] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K.Q. (eds.) *Advances in Neural Information Processing Systems*, vol. 27. Curran Associates, Inc., ??? (2014)
- [115] Zhu, J.-Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* (2017)
- [116] Park, T., Efros, A.A., Zhang, R., Zhu, J.-Y.: Contrastive learning for unpaired image-to-image translation. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX* 16, pp. 319–345 (2020). Springer
- [117] Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020)
- [118] Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* **34**, 8780–8794 (2021)
- [119] Chen, J., He, Y., Frey, E.C., Li, Y., Du, Y.: Vit-v-net: Vision transformer for unsupervised volumetric medical image registration. *arXiv preprint*

- [120] Zhang, Y., Pei, Y., Zha, H.: Learning dual transformer network for diffeomorphic registration. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part IV 24, pp. 129–138 (2021). Springer
- [121] Mok, T.C., Chung, A.: Affine medical image registration with coarse-to-fine vision transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20835–20844 (2022)
- [122] Chen, J., Frey, E.C., He, Y., Segars, W.P., Li, Y., Du, Y.: Transmorph: Transformer for unsupervised medical image registration. *Medical image analysis* **82**, 102615 (2022)
- [123] Song, L., Liu, G., Ma, M.: Td-net: unsupervised medical image registration network based on transformer and cnn. *Applied Intelligence* **52**(15), 18201–18209 (2022)
- [124] Wang, Y., Qian, W., Li, M., Zhang, X.: A transformer-based network for deformable medical image registration. In: Artificial Intelligence: Second CAAI International Conference, CICA 2022, Beijing, China, August 27–28, 2022, Revised Selected Papers, Part I, pp. 502–513 (2022). Springer
- [125] Shi, J., He, Y., Kong, Y., Coatrieux, J.-L., Shu, H., Yang, G., Li, S.: Xmorpher: Full transformer for deformable medical image registration via cross attention. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 217–226 (2022). Springer
- [126] Zhu, Y., Lu, S.: Swin-voxlomorph: A symmetric unsupervised learning model for deformable medical image registration using swin transformer. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 78–87 (2022). Springer
- [127] Chen, Z., Zheng, Y., Gee, J.C.: Transmatch: A transformer-based multi-level dual-stream feature matching network for unsupervised deformable image registration. *IEEE Transactions on Medical Imaging* (2023)
- [128] Wang, H., Ni, D., Wang, Y.: Modet: Learning deformable image registration via motion decomposition transformer. *arXiv preprint arXiv:2306.05688* (2023)
- [129] Kamran, S.A., Fariha Hossain, K., Tavakkoli, A., Zuckerbrod, S., Baker, S.A., Sanders, K.M.: Fundus2angio: a conditional gan architecture for generating fluorescein angiography images from retinal fundus photography. In: Advances in Visual Computing: 15th International Symposium, ISVC 2020, San Diego, CA, USA, October 5–7, 2020, Proceedings, Part II 15, pp. 125–138 (2020). Springer

- [130] Andreini, P., Ciano, G., Bonechi, S., Graziani, C., Lachi, V., Mecocci, A., Sodi, A., Scarselli, F., Bianchini, M.: A two-stage gan for high-resolution retinal image generation and segmentation. *Electronics* **11**(1), 60 (2021)
- [131] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: *International Conference on Learning Representations* (2021)
- [132] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022 (2021)
- [133] Hering, A., Hansen, L., Mok, T.C.W., Chung, A.C.S., Siebert, H., Häger, S., Lange, A., Kuckertz, S., Heldmann, S., Shao, W., Vesal, S., Rusu, M., Sonn, G., Estienne, T., Vakalopoulou, M., Han, L., Huang, Y., Yap, P.-T., Brudfors, M., Balbastre, Y., Joutard, S., Modat, M., Lifshitz, G., Raviv, D., Lv, J., Li, Q., Jaouen, V., Visvikis, D., Fourcade, C., Rubeaux, M., Pan, W., Xu, Z., Jian, B., De Benetti, F., Wodzinski, M., Gunnarsson, N., Sjölund, J., Grzech, D., Qiu, H., Li, Z., Thorley, A., Duan, J., Großbröhm, C., Hoopes, A., Reinertsen, I., Xiao, Y., Landman, B., Huo, Y., Murphy, K., Lessmann, N., Van Ginneken, B., Dalca, A.V., Heinrich, M.P.: Learn2reg: comprehensive multi-task medical image registration challenge, dataset and evaluation in the era of deep learning. *IEEE Transactions on Medical Imaging*, 1–1 (2022) <https://doi.org/10.1109/TMI.2022.3213983>
- [134] Ding, L., Kang, T.D., Kuriyan, A.E., Ramchandran, R.S., Wykoff, C.C., Sharma, G.: Combining feature correspondence with parametric chamfer alignment: Hybrid two-stage registration for ultra-widefield retinal images. *IEEE Transactions on Biomedical Engineering* **70**(2), 523–532 (2023) <https://doi.org/10.1109/TBME.2022.3196458>
- [135] Wykoff, C.C., Nittala, M.G., Zhou, B., Fan, W., Velaga, S.B., Lampen, S.I.R., Rusakevich, A.M., Ehlers, J.P., Babiuch, A., Brown, D.M., Ip, M.S., Sadda, S.R., Wykoff, C.C., Nittala, M.G., Zhou, B., Fan, W., Velaga, S.B., Rusakevich, A.M., Lampen, S.I.R., Ip, M.S., Sadda, S.R., Ehlers, J.P., Srivastava, S.K., Reese, J.L., Babiuch, A., Talcott, K., Figueiredo, N., Hach, J., Ou, W.C., Fish, R.H., Benz, M.S., Chen, E., Kim, R.Y., Major, J.C., O'Malley, R.E., Brown, D.M., Shah, A.R., Scheffler, A.C., Wong, T.P., Henry, C.R.: Intravitreal aflibercept for retinal nonperfusion in proliferative diabetic retinopathy: Outcomes from the randomized recovery trial. *Ophthalmology Retina* **3**(12), 1076–1086 (2019) <https://doi.org/10.1016/j.oret.2019.07.011>