

Multiview 2D/3D Rigid Registration via a Point-Of-Interest Network for Tracking and Triangulation

Haofu Liao^{1*}, Wei-An Lin^{2*}, Jiarui Zhang³, Jingdan Zhang⁴, Jiebo Luo¹, S. Kevin Zhou^{5,6}

¹University of Rochester

²University of Maryland, College Park

³Rutgers University

⁴Z2W Corporation

⁵Chinese Academy of Sciences

⁶Peng Cheng Laboratory, Shenzhen

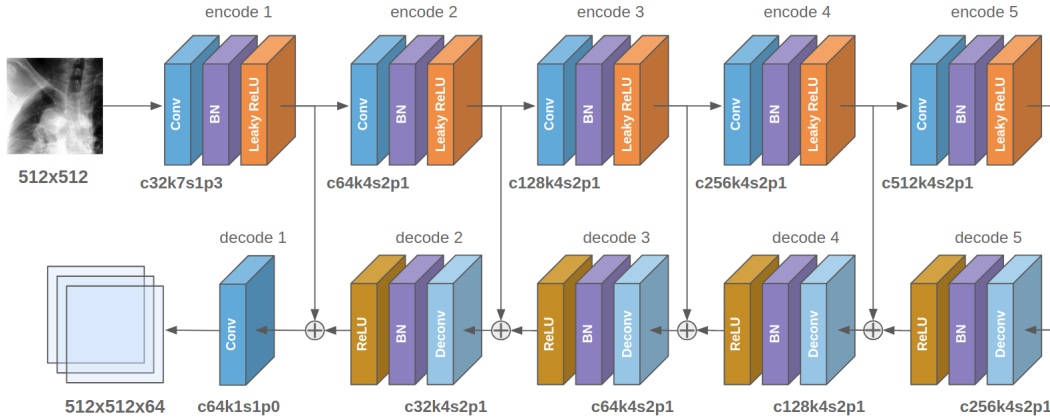


Figure 1: The detailed architecture of the U-Net branch of the POINT network.

A. Network Architecture

Figure 1 shows details of the U-Net branches of the POINT network. There are five pairs of encoding and decoding blocks in the network. Each pair is connected with a skip-connection that shuttles the outputs of the encoding block to the decoding block directly. The layer configurations are given below the convolutional and deconvolutional layers, where “c”, “k”, “s” and “p” denote the channel, kernel, stride and padding size, respectively. The “+” symbol at the end of a skip-connection denotes a concatenation operation that combines feature maps from the encoding block “encode n ” and a lower level decoding block “decode $n+1$ ”.

Figure 2 illustrates the feature extraction process for a 3×3 feature kernel. Given a feature map \mathbf{F}^D of size $M \times N \times C$ obtained from the DRR U-Net branch and a set of POIs $\mathcal{P} = \{\mathbf{x}_1^D, \mathbf{x}_2^D, \dots, \mathbf{x}_m^D\}$, the feature extraction layer first interpolates each channel of the feature map \mathbf{F}^D at $\forall \mathbf{x}_i^D \in \mathcal{P}$, which gives a feature vector of size $1 \times 1 \times C$. Then, to obtain a richer feature representation, it also interpolates the neighborhood of \mathbf{x}_i^D which results a feature kernel $\mathbf{F}^D(\mathbf{x}_i^D)$ of size $3 \times 3 \times C$.

* indicates equal contributions.

B. Derivation

In this section, we show how

$$\mathbf{D}(\mathbf{x})\mathbf{R}_{\text{view}}\mathbf{X} = \mathbf{c}\mathbf{x} - \mathbf{D}(\mathbf{x})\mathbf{t}_{\text{view}} \quad (1)$$

can be derived from

$$\mathbf{x}' = \mathbf{K} [\mathbf{R}_{\text{view}} \quad \mathbf{t}_{\text{view}} + \mathbf{h}] \begin{pmatrix} \mathbf{X} \\ 1 \end{pmatrix}. \quad (2)$$

Proof. Let

$$\tilde{\mathbf{X}} = \mathbf{R}_{\text{view}}\mathbf{X} + \mathbf{t}_{\text{view}} \quad (3)$$

and recall that

$$\mathbf{K} = \begin{bmatrix} -d\mathbf{I} & 0 \\ \mathbf{0} & 1 \end{bmatrix}, \quad (4)$$

then (2) can be written as

$$\begin{aligned} \mathbf{x}' &= \begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \begin{bmatrix} -d\mathbf{I} & 0 \\ \mathbf{0} & 1 \end{bmatrix} (\tilde{\mathbf{X}} + \mathbf{h}) \\ &= \begin{bmatrix} [-d\mathbf{I} & 0] (\tilde{\mathbf{X}} + \mathbf{h}) \\ [\mathbf{0} & 1] (\tilde{\mathbf{X}} + \mathbf{h}) \end{bmatrix} \end{aligned} \quad (5)$$

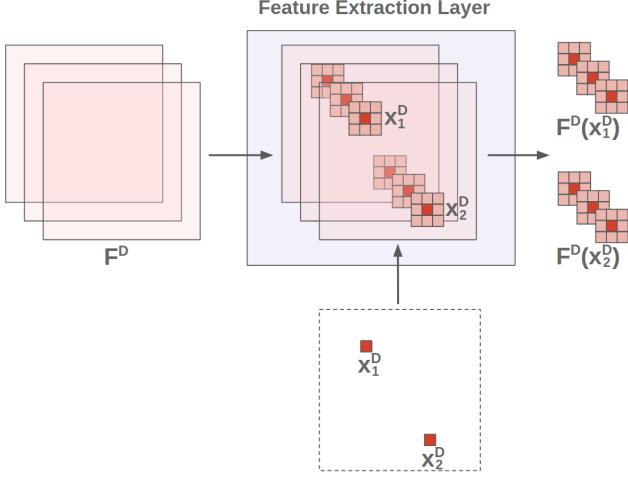


Figure 2: An illustration of the feature extraction process of a 3×3 kernel.

Since $\mathbf{x} = (x'/z', y'/z')^T$, we can get

$$\mathbf{x} = \begin{bmatrix} -d\mathbf{I} & 0 \\ \mathbf{0} & 1 \end{bmatrix} (\tilde{\mathbf{X}} + \mathbf{h}) \quad (6)$$

$$[d\mathbf{I} \quad \mathbf{x}] \tilde{\mathbf{X}} = [-d\mathbf{I} \quad -\mathbf{x}] \mathbf{h} \quad (7)$$

Also, recall that

$$\mathbf{h} = \begin{pmatrix} 0 \\ 0 \\ -c \end{pmatrix}, \mathbf{D}(\mathbf{x}) = \begin{bmatrix} d & 0 & \mathbf{x} \\ 0 & d & \mathbf{x} \end{bmatrix}, \quad (8)$$

we can simplify (7) as

$$\mathbf{D}(\mathbf{x})\tilde{\mathbf{X}} = c\mathbf{x}. \quad (9)$$

Finally, using (3), we have

$$\mathbf{D}(\mathbf{x})\mathbf{R}_{\text{view}}\mathbf{X} = c\mathbf{x} - \mathbf{D}(\mathbf{x})\mathbf{t}_{\text{view}} \quad (10)$$

□

C. POI Tracking Results

See next page.

