

# Collaborative Quantization Embeddings for Intra-Subject Prostate MR Image Registration

Ziyi Shen<sup>1</sup>, Qianye Yang<sup>1</sup>, Yuming Shen<sup>2</sup>, Francesco Giganti<sup>3</sup>, Vasilis Stavrinides<sup>1</sup>, Richard Fan<sup>4</sup>, Caroline Moore<sup>1</sup>, Mirabela Rusu<sup>4</sup>, Geoffrey Sonn<sup>4</sup>, Philip Torr<sup>2</sup>, Dean Barratt<sup>1</sup>, and Yipeng Hu<sup>1</sup>

<sup>1</sup> University College London, London, UK

<sup>2</sup> University of Oxford, Oxford, UK

<sup>3</sup> University College London Hospital NHS Foundation Trust, London, UK

<sup>4</sup> Stanford University, CA 94305, USA

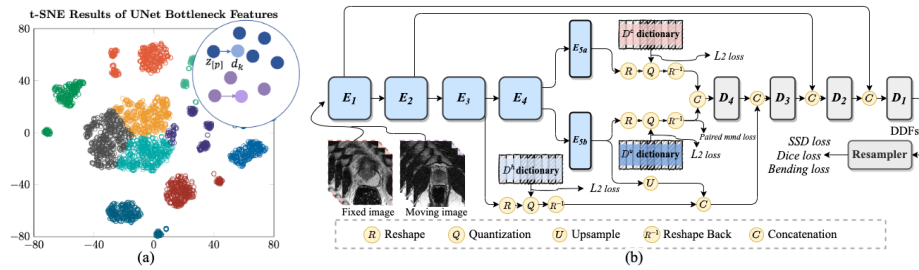
joanshen0508@gmail.com

**Abstract.** Image registration is useful for quantifying morphological changes in longitudinal MR images from prostate cancer patients. This paper describes a development in improving the learning-based registration algorithms, for this challenging clinical application often with highly variable yet limited training data. First, we report that the latent space can be clustered into a much lower dimensional space than that commonly found as bottleneck features at the deep layer of a trained registration network. Based on this observation, we propose a hierarchical quantization method, discretizing the learned feature vectors using a jointly-trained dictionary with a constrained size, in order to improve the generalisation of the registration networks. Furthermore, a novel collaborative dictionary is independently optimised to incorporate additional prior information, such as the segmentation of the gland or other regions of interest, in the latent quantized space. Based on 216 real clinical images from 86 prostate cancer patients, we show the efficacy of both the designed components. Improved registration accuracy was obtained with statistical significance, in terms of both Dice on gland and target registration error on corresponding landmarks, the latter of which achieved 5.46 mm, an improvement of 28.7% from the baseline without quantization. Experimental results also show that the difference in performance was indeed minimised between training and testing data.

**Keywords:** Registration · Quantization · Prostate Cancer.

## 1 Introduction

Whilst the diagnostic value in multiparametric MR imaging for prostate cancer, before or after biopsy for histopathology examination, has been identified [14,18], attention has been quickly turned to using this non-invasive imaging technique to monitor the disease at its early stage. Many have speculated that the temporal changes in morphology and intensity pattern in prostate gland can indicate the progression of the cancer [7]. Establishing spatial correspondence between two



**Fig. 1.** Demonstration of the proposed prostate registration framework. (a) The t-SNE [10] visualization of the encoder outputs of a U-Net-like prostate registration network. (b) The network structure of HiCo-Net.

or more images, medical image registration has been proposed for aligning MR images from prostate cancer patients acquired at different time points. The estimated intra-subject spatial transformation, representing corresponding spatial locations between prostate glands, is an important quantitative tool to track the radiological evolution of prostate glands [24]. Aligning anatomical structures in lower pelvic region has been recognised to be challenging using ‘classical’ iterative algorithms [27], perhaps due to the highly patient-specific imaging characteristics in organs, such as prostate glands being having distinct patient-specific intensity patterns on T2-weighted MR images. Recent development has focused on learning-based algorithms for their effective and efficient inference [5,1,12,21,6,22]. Empirically, only a few ‘types’ of MR image features are reliably useful for establishing correspondence, such as volume, shape, anatomical and pathological structures, within or around the prostate gland and potential cancerous regions [24]. In addition, prostate capsules and different types of pathology are known to be highly variable and specific to individual patients [19,28]. Therefore, features are more likely to be ‘easy-to-learn’ to the distinct inter-subject differences.

We argue that these two above intuitions may warrant a compact feature adequately containing intra-subject correspondence for this application, although a deep network may still be required to learn such a representation [8]. However, deep models are over-parametrized [13], where the hidden representation may still carry information that is not related the task. With limited MR training data in a real clinical application, this redundancy and over-parametrization degrade the learned features’ ability to generalize, leading to overfitting.

To illustrate this redundancy, we visualize the t-SNE [10] results of the deep features for registering prostate MR images in Fig. 1 (a). Given an  $N$ -sample prostate image set, we extract the bottleneck features of all samples, producing  $N$  feature maps with a size of  $W \times H \times T \times C$ . Here,  $W \times H \times T$  refers to the shape of the output feature, while  $C$  is the channel number. Each  $C$ -dimensional vector represents a super-pixel in the encoded feature. We visualize the t-SNE result of

all super-pixels in the set. It is shown in Fig. 1 (a) that the features are scattered into a limited numbers of groups. In other words, one can roughly represent the whole feature space using a smaller number of latent topics, consistent with the above discussion in limited corresponding structures and inter-subject variability.

This observation opens the door to compress the features with deep vector quantization (VQ) [15,16,20]. We refer to the anatomical knowledge of prostate, whose appearance varies from subjects, and propose to represent the neural features with a small set of vocabulary vectors, where useful anatomical and pathological structures are preserved. This ideology has been proved effective with limited training samples [4], which is indeed the case our task.

In this paper, the aforementioned motivation drives us to a VQ-based prostate registration solution. Specifically, we apply VQ to the middle of a registration U-Net [17] to make effective use of the feature space. The dictionary of a vector quantizer is usually data-driven. To efficiently improve the anatomical and pathological awareness of the model, we introduce another quantizer in parallel to the randomly initialized one, of which the dictionary preserves the specific local features of prostate interior. This is done by abstracting knowledge from a deep prostate segmentation network. The combination of the data-driven and the pre-defined quantizers is termed **collaborative quantization**. In addition, we explore the multi-scale structure of prostate MR images to fit both the global and local patterns of a moving image to the fixed one, which suggests a **hierarchical quantization** structure similar to [16]. Therefore, we name our model as Hierarchically & Collaboratively Quantized Network (HiCo-Net).

Our contributions can be summarised as follows: (1) We propose a feature quantization framework as regularization to alleviate the gap between training and test data for registration (2) We introduce a collaborative quantizer that encodes structure features of the gland boundary to better represent lesions and landmarks. (3) Our experiments show that the proposed HiCo-Netsuccessfully relieves the overfitting problem in weakly-supervised longitudinal prostate image registration, and outperforms the state-of-the-art.

## 2 Method

We consider a pairwise MR image registration problem. Let  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$  be the collection of images pairs of prostate, where  $\mathbf{x}_i = (\mathbf{x}_i^s, \mathbf{x}_i^t)$  denotes a pair, with  $n$  being the number of image pairs.  $\mathbf{x}_i^s$  and  $\mathbf{x}_i^t$  respectively refers to the moving image and the fixed one. Each image pair comes with a pair of prostate gland anatomical segmentation maps  $(\mathbf{m}_i^s, \mathbf{m}_i^t)$  for weakly-supervised training. For each  $\mathbf{x}_i$ , the goal is to predict a dense displacement field (DDF)  $\mathbf{u}_i$  to establish voxel-level correspondence.

### 2.1 Preliminary: Deep Vector Quantization

VQ [20] quantizes an arbitrary representation tensor using a fixed number of values defined by a dictionary  $\mathcal{D} = \{\mathbf{d}_i\}_{i=1}^K$ ,  $\mathbf{d}_i \in \mathbb{R}^C$ , where  $K$  is the dictionary size and  $C$  is the dimensionality of each code. Specifically, an output of

an encoder  $E(\mathbf{x}) \in \mathbb{R}^{H \times W \times T \times C}$  is obtained by passing an MR image pair  $\mathbf{x}$  through a CNN. We are going to quantize each  $C$ -dimensional vector of  $E(\mathbf{x})$ . For simplicity, the rest of this paper denotes a voxel super-position in a raster scan order with a coordinate  $p = 1 \cdots HWT$ , i.e.,  $E(\mathbf{x})_{[p]} \in \mathbb{R}^C$ . Hence, a vector quantization operator  $Q(\cdot)$  can be defined as follows:

$$\mathbf{z}_{[p]} = Q(E(\mathbf{x})_{[p]}; \mathcal{D}) = \mathbf{d}_k, \text{ where } k = \underset{i}{\operatorname{argmin}} \|E(\mathbf{x})_{[p]} - \mathbf{d}_i\|. \quad (1)$$

Then,  $\mathbf{z}$  replaces  $E(\mathbf{x})$  and is forwarded to the rest of the network. The encoder receives the gradients from top of the quantizer through the straight-through estimator [2], i.e.,  $\partial \mathbf{z} / \partial E := \mathbb{I}$ . VQ incorporates two additional loss terms to enforce the output of the encoder to be similar to the quantized results:

$$\mathcal{L}_Q(E(\mathbf{x}), \mathcal{D}) = \sum_p \|\operatorname{sg}(E(\mathbf{x})_{[p]}) - \mathbf{z}_{[p]}\|_2^2 + \beta \|E(\mathbf{x})_{[p]} - \operatorname{sg}(\mathbf{z}_{[p]})\|_2^2, \quad (2)$$

where  $\operatorname{sg}(\cdot)$  is the stop-gradient operator and  $\beta = 0.25$  is the hyperparameter.

## 2.2 Model Overview

Fig. 1 (b) depicts the schematic of HiCo-Net. It generally undergoes a U-Net-like structure with an encoder  $E(\cdot)$  and a decoder  $D(\cdot)$ . An image pair  $\mathbf{x} = (\mathbf{x}^s, \mathbf{x}^t)$  is firstly concatenated together and then rendered to the encoder, while the decoder produces  $\mathbf{u}$ , the desired DDF. We particularly denote the output of each residual block as  $E_1(\mathbf{x}), E_2(\mathbf{x}) \cdots$ . Notably,  $E_4(\mathbf{x})$  is rendered to two parallel convolutional layers,  $E_{5a}(\mathbf{x})$  and  $E_{5b}(\mathbf{x})$ , that are followed by two quantizers. We respectively term them the collaborative quantizer (Sec. 2.5) and the vanilla one (Sec. 2.3). The skip connection between  $E_3(\mathbf{x})$  and the decoder is quantized as well. Since it also mixes multi-scale information from  $E_{5b}(\mathbf{x})$  afterwards, we name it the hierarchical quantizer (Sec. 2.4). The intuition behind this design is given in their respective sections as follows. The DDF output then contributes to the conventional weakly-supervised prostate registration losses with a resampler.

## 2.3 Vanilla Quantization

We first introduce a vanilla quantizer that quantizes the output of  $E_{5b}(\mathbf{x})$ . It behaves identical to the original VQ operation [20]. Shown in Fig. 1 (b), its dictionary  $\mathcal{D}^v$  is randomly initialized and is updated by back-propagation during training. In this way, the global information, a relatively fixed structure of the MR prostate images, is regularized for better generalization to test data. We denote the quantization loss for the vanilla quantizer as  $\mathcal{L}_V(\mathbf{x}) = \mathcal{L}_Q(E_{5b}(\mathbf{x}), \mathcal{D}^v)$ .

## 2.4 Hierarchical Quantization

Image features of a deep network often carry local information, which can benefit from multi-scale modelling for positional alignment. A hierarchical representation quantizer has been proved to be effective to perceiving this [16].

To implement this, we employ a hierarchical quantizer to quantize the output of  $E_3(\mathbf{x})$ , of which the dictionary is denoted as  $\mathcal{D}^h$ . The quantized result is added by the output of  $E_4(\mathbf{x})$ . Since the voxel sizes of them mismatch, one needs to firstly upsample  $E_4(\mathbf{x})$ , as is shown in Fig. 1 (b).  $\mathcal{D}^h$  is randomly initialized. The hierarchical quantizer introduces a quantization loss as  $\mathcal{L}_H(\mathbf{x}) = \mathcal{L}_Q(E_3(\mathbf{x}), \mathcal{D}^h)$ .

## 2.5 Collaborative Quantization

As is discussed in Sec. 1, the awareness of prostate contour is of key importance to this prostate registration task. This inspires us to transfer knowledge to our model from a deep segmentation network, without requiring segmentation during inference. VQ allows us to conveniently initialize the dictionary values according to the segmentation network’s output as prior knowledge. In particular, we first train a U-Net-based segmentation network on the training dataset, and then extract an  $H \times W \times T \times C$  tensor for each image with its encoder. Each  $C$ -dimensional vector of all images’ features is treated as an instance for K-means clustering. We cache the values of  $K$  cluster centers produced by K-means, and use them to initialize the dictionary of the collaborative quantizer, i.e.,  $\mathcal{D}^c$ . An analogy of this procedure is provided in the supplemental.

The collaborative quantizer takes input from  $E_{5a}(\mathbf{x})$  as a compensation to the fully-data-driven vanilla quantizer by concatenating their outputs afterwards. We similarly define its quantization loss as  $\mathcal{L}_C(\mathbf{x}) = \mathcal{L}_Q(E_{5a}(\mathbf{x}), \mathcal{D}^c)$ .

## 2.6 Training

**Quantization Loss.** The training objective of HiCo-Net is a combination of the quantization losses and the conventional weakly-supervised registration ones. We first define the overall quantization loss of an image pair  $\mathbf{x}$ :

$$\mathcal{L}_{\text{Quant}}(\mathbf{x}) = \mathcal{L}_V(\mathbf{x}) + \mathcal{L}_C(\mathbf{x}) + \mathcal{L}_H(\mathbf{x}). \quad (3)$$

The three quantization terms above can have equal weights as they are mutually independent, and are imposed to different stages of the registration network.

**SSD Loss.** The Sum-of-Square Differences (SSD) loss [1] measures the similarity between the translated image and the fixed one. One needs to firstly resample  $\mathbf{x}^s$  using the DDF  $\mathbf{u}$  and then compute

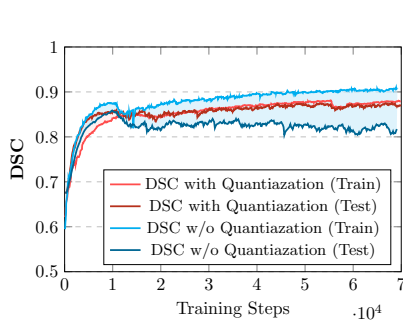
$$\mathcal{L}_{\text{SSD}}(\mathbf{x}) = \|\mathbf{u} \otimes \mathbf{x}^s - \mathbf{x}^t\|_2^2, \quad (4)$$

where  $\otimes$  refers to the resampling operation.

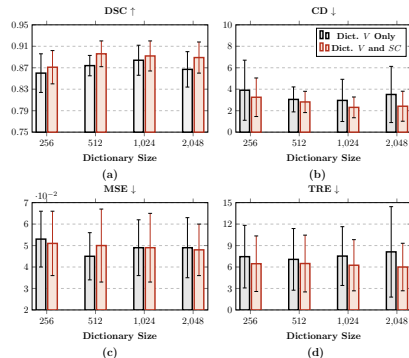
**Dice Loss.** This loss has shown effectiveness in aligning organ shapes and positions [23], and is applied to the masks:

$$\mathcal{L}_{\text{Dice}}(\mathbf{x}) = -\text{Dice}(\mathbf{u} \otimes \mathbf{m}^s, \mathbf{m}^t). \quad (5)$$

**Bending Regularization.** We use this regularization term  $\mathcal{L}_{\text{Bend}}(\mathbf{u})$  [26] to penalise the non-smoothness of the generated DDF.



**Fig. 2.** Illustration of how  $\mathcal{D}^c$  is initialized before training.



**Fig. 3.** Ablation study of dictionary size.

**Overall Training Objective.** By defining the losses above, we can simply compose a linear combination of them as the final loss of HiCo-Net as follows:

$$\mathcal{L}_{\text{All}}(\mathbf{x}) = \lambda_Q \mathcal{L}_{\text{Quant}}(\mathbf{x}) + \lambda_S \mathcal{L}_{\text{SSD}}(\mathbf{x}) + \lambda_D \mathcal{L}_{\text{Dice}}(\mathbf{x}) + \lambda_B \mathcal{L}_{\text{bend}}(\mathbf{u}), \quad (6)$$

where  $\lambda_Q = 1$ ,  $\lambda_S = 1$ ,  $\lambda_D = 1$  and  $\lambda_B = 50$  are hyperparameters. Our models are trained with stochastic gradient descent algorithms.

## 3 Experiment

### 3.1 Experimental Settings

**Implementation.** We use a basic U-Net equipped with skip connections between encoder and decoder. Our encoder consists of 4 residual blocks, a total of 12 convolutional layers. In addition, we add 2 convolutional layers to expand heads for the subsequent hierarchical and collaborative quantization operators as shown in Fig. 1 (b). The channel size of the hierarchical dictionary  $\mathcal{D}^h$  is 128. The vanilla dictionary  $\mathcal{D}^v$  and collaborative one  $\mathcal{D}^c$  both have a feature channel number of 256. We by default set the vocabulary size of  $\mathcal{D}^h$  and  $\mathcal{D}^v$  to 1024, while fixing the size of  $\mathcal{D}^c$  one to 512. As per the initialization of  $\mathcal{D}^c$ , we use a U-Net as our segmentation network, which is utilized to collect features for K-means clustering. Note that the segmentation network is not involved in our registration training. We set training batch size to 4, and use the Adam optimizer with a learning rate of  $10^{-4}$ . The network is trained for 1000 epochs at most, taking three days on an NVIDIA Tesla V100 GPU. The network architecture and code are available at <https://github.com/joanshen0508/HiCo-Net>.

**Dataset.** The utilized dataset consists of 216 longitudinal prostate T2-weighted MR images from 86 patients, acquired from University College London Hospitals NHS Foundation. It is divided into three folds, containing 70, 6, and 10 patients for training, validation, and test. Each patient has 2-4 images, with

**Table 1.** Ablation study of hierarchical and collaborative quantization.

$\mathcal{D}^v$	$\mathcal{D}^h$	$\mathcal{D}^c$	DSC	CD	MSE	TRE
w/o registration			0.700±0.097	12.63±5.810	0.051±0.014	13.72±5.833
			0.859±0.038	4.187±2.050	0.049±0.013	7.657±4.212
✓			0.884±0.028	2.958±1.967	0.049±0.013	7.529±4.109
✓			0.887±0.264	2.644±1.469	0.048±0.014	6.158±3.539
✓			0.865±0.027	3.011±1.635	0.050±0.015	7.551±4.435
✓ w/o pretrain			0.892±0.028	2.308±0.967	0.049±0.016	6.248±3.577
✓			0.881±0.025	3.091±1.557	0.043±0.013	5.457±3.489

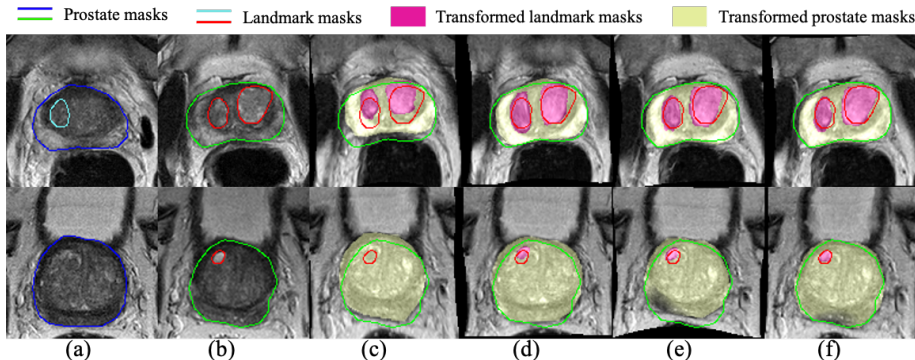
an average interval between consecutive visits of 18.1 and a standard deviation of 10.3 months. Before training, we resample the data to  $0.7 \times 0.7 \times 0.7mm^3$  and normalize the intensity to  $[0, 1]$ . To train the proposed prostate registration model, we also crop the dataset and generate the dataset with the size of  $128 \times 128 \times 102$ . On the test set, 141 anatomical and pathological landmarks are manually identified on moving and fixed images, including patient-specific fluid-filled cysts, calcification and centroids of zonal boundaries.

**Evaluation Metrics.** We adopt the conventional weakly-supervised registration metrics including Dice Similarity Coefficient (DSC) and Centroid Distance (CD) between the prostate glands. The Mean-Squared Error (MSE) between the fixed image and wrapped moving image is as well reported. Registration should support downstream clinical image analysis task. To demonstrate the effectiveness of our method, we further report the Target Registration Error (TRE), which calculates the difference of landmarks between fixed image and predicted result.

### 3.2 Ablation study

**The Effect of Feature Quantization.** We first build a registration model without quantization similar to [24] and compare it with a variant of HiCo-Net that only involves the vanilla quantizer. As shown in Fig. 2, our quantized version effectively narrows the accuracy gap between training and test, observing no overfitting problem. We also report the performance of HiCo-Net with different combinations of the three quantizers in Tab. 1. Compared with the unquantized baseline, the TRE is reduced from  $7.657 \pm 4.212$  mm to  $6.248 \pm 3.577$  mm ( $p\_value = 0.0001$  under paired t-test) when applying the vanilla and collaborative quantization, and it further decreases to  $5.457 \pm 3.489$  mm ( $p\_value < 0.0001$ ) when employing all the three quantizers.

**Hierarchical and Collaborative Quantization.** The hierarchical quantization scheme mixes the global and local information, and obtains the best results in Tab. 1. We also consider randomly initializing the values of  $\mathcal{D}^c$ . Its gain against the single-quantizer baseline is marginal, but once initialized by segmentation feature vectors, the collaborative embedding improves the spatial alignment to focus on the local semantic discrepancy, and obtains a better registration performance. We provide qualitative comparison results in Fig. 4. The proposed



**Fig. 4. Effect of proposed method.** (a) moving image. (b) fixed image. (c) w/o quantization. (d) w/  $\mathcal{D}^v$ . (e) w/  $\mathcal{D}^v$  and  $\mathcal{D}^h$ . (f) w/  $\mathcal{D}^v$ ,  $\mathcal{D}^h$  and  $\mathcal{D}^c$ .

**Table 2.** Comparison with the-state-of-the-arts prostate registration methods.

Method	DSC	CD	MSE	TRE	Run-time
NiftyReg [11]	0.270±0.304	22.869±11.761	0.041±0.019	21.147±15.841	45.76
VoxelMorph[1]	0.763±0.081	8.842±3.156	0.053±0.015	8.833±5.147	0.69
DeepTag[25]	0.822±0.083	7.594±2.905	0.052±0.013	7.458±4.815	1.95
Contrastive[9]	0.856±0.117	4.973±2.407	0.054±0.018	8.2166±4.407	0.31
Basic U-Net	0.859±0.038	4.187±2.050	0.049±0.013	7.657±4.212	0.62
VAE-like	0.865±0.029	3.623±2.189	0.045±0.019	7.626±3.948	0.72
<b>HiCo-Net</b>	0.881±0.025	3.091±1.557	0.043±0.013	5.457±3.489	0.68

method performs well on aligning local patterns to the fixed image. We evaluate the dictionary size of  $\mathcal{D}^v$  and  $\mathcal{D}^c$ , shown in Fig. 3, which suggests a dictionary size of 512 and 1024 for  $\mathcal{D}^v$  and  $\mathcal{D}^c$  respectively is the best option.

**Inter-Subject Extension.** To explore quantization for further generic application, we also validate the proposed method on inter-subject prostate MR data [3]. We notice that the performance increases with quantization (DSC:  $0.80 \pm 0.11 \rightarrow 0.86 \pm 0.04$ , CD:  $4.17 \pm 2.43 \rightarrow 2.12 \pm 1.33$ , MSE:  $0.04 \pm 0.02 \rightarrow 0.03 \pm 0.02$ ). This task is challenging as the presence of prostate varies from different identities.

### 3.3 Comparison with Existing Methods

We compare HiCo-Net with a non-optimised iterative method [11] and some well-known deep registration methods [1,25,9]. To further validate the encoder-decoder structure, a common U-net and a VAE framework are implemented for prostate registration. As shown in Tab. 2, the proposed method obtains competitive results in all metrics. Remarkably, the number of negative Jacobian determinants of our method is  $0.0 \pm 0.0$ . The consuming time is also reported. In



addition, our collaborative quantization algorithm is free from additional sub-network embedding, avoiding large memory consumption.

## 4 Conclusion

In this paper, we proposed a collaborative quantization framework for prostate MR image registration, which was named HiCo-Net. We introduced a hierarchical quantizer that jointly regularizes the global and local latent information to benefit the displacement prediction. In addition, we designed a collaborative dictionary that was equipped with helpful anatomical structure knowledge to perceive the local semantic discrepancy. The experiments showed that this method performed favorably against state-of-the-art registration methods and bypassed the overfitting problem for our dataset with a moderate size. Representing and quantizing inter-subject cues for registration can be our future work.

**Acknowledgements** This work was supported by the International Alliance for Cancer Early Detection, an alliance between Cancer Research UK [C28070/A30912; C73666/A31378], Canary Center at Stanford University, the University of Cambridge, OHSU Knight Cancer Institute, University College London and the University of Manchester. This work was also supported by the Wellcome/EPSRC Centre for Interventional and Surgical Sciences [203145Z/16/Z].

## References

1. Balakrishnan, G., Zhao, A., Sabuncu, M.R., Guttag, J., Dalca, A.V.: Voxelmorph: a learning framework for deformable medical image registration. *IEEE transactions on medical imaging* **38**(8), 1788–1800 (2019) [2](#), [5](#), [8](#)
2. Bengio, Y., Léonard, N., Courville, A.: Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432* (2013) [4](#)
3. Bloch, N., Madabhushi, A., Huisman, H., Freymann, J., Kirby, J., Grauer, M., Enquobahrie, A., Jaffe, C., Clarke, L., Farahani, K.: Nci-isbi 2013 challenge: automated segmentation of prostate structures. *The Cancer Imaging Archive* **370**(6), 5 (2015) [8](#)
4. Chen, K., Lee, C.G.: Incremental few-shot learning via vector quantization in deep embedded space. In: *ICLR* (2021) [3](#)
5. Chen, X., Meng, Y., Zhao, Y., Williams, R., Vallabhaneni, S.R., Zheng, Y.: Learning unsupervised parameter-specific affine transformation for medical images registration. In: *MICCAI*. pp. 24–34. Springer (2021) [2](#)
6. 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) [2](#)
7. Kim, C.K., Park, B.K., Lee, H.M., Kim, S.S., Kim, E.: Mri techniques for prediction of local tumor progression after high-intensity focused ultrasonic ablation of prostate cancer. *American Journal of Roentgenology* **190**(5), 1180–1186 (2008) [1](#)

8. Liu, F., Yan, K., Harrison, A.P., Guo, D., Lu, L., Yuille, A.L., Huang, L., Xie, G., Xiao, J., Ye, X., et al.: Same: Deformable image registration based on self-supervised anatomical embeddings. In: MICCAI. pp. 87–97. Springer (2021) [2](#)
9. Liu, L., Aviles-Rivero, A.I., Schönlieb, C.B.: Contrastive registration for unsupervised medical image segmentation. arXiv preprint arXiv:2011.08894 (2020) [8](#)
10. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. *Journal of Machine Learning Research* **9**(Nov), 2579–2605 (2008) [2](#)
11. Modat, M., Ridgway, G.R., Taylor, Z.A., Lehmann, M., Barnes, J., Hawkes, D.J., Fox, N.C., Ourselin, S.: Fast free-form deformation using graphics processing units. *Computer methods and programs in biomedicine* **98**(3), 278–284 (2010) [8](#)
12. Mok, T.C., Chung, A.: Large deformation diffeomorphic image registration with laplacian pyramid networks. In: MICCAI. pp. 211–221. Springer (2020) [2](#)
13. Molchanov, D., Ashukha, A., Vetrov, D.: Variational dropout sparsifies deep neural networks. In: ICML. pp. 2498–2507. PMLR (2017) [2](#)
14. Moore, C.M., Giganti, F., Albertsen, P., Allen, C., Bangma, C., Briganti, A., Carroll, P., Haider, M., Kasivisvanathan, V., Kirkham, A., et al.: Reporting magnetic resonance imaging in men on active surveillance for prostate cancer: the precise recommendations—a report of a european school of oncology task force. *European urology* **71**(4), 648–655 (2017) [1](#)
15. Peng, J., Liu, D., Xu, S., Li, H.: Generating diverse structure for image inpainting with hierarchical vq-vae. In: CVPR. pp. 10775–10784 (2021) [3](#)
16. Razavi, A., Van den Oord, A., Vinyals, O.: Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems* **32** (2019) [3](#), [4](#)
17. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015) [3](#)
18. Schoots, I.G., Petrides, N., Giganti, F., Bokhorst, L.P., Rannikko, A., Klotz, L., Villers, A., Hugosson, J., Moore, C.M.: Magnetic resonance imaging in active surveillance of prostate cancer: a systematic review. *European urology* **67**(4), 627–636 (2015) [1](#)
19. Song, X., Guo, H., Xu, X., Chao, H., Xu, S., Turkbey, B., Wood, B.J., Wang, G., Yan, P.: Cross-modal attention for mri and ultrasound volume registration. In: MICCAI. pp. 66–75. Springer (2021) [2](#)
20. Van Den Oord, A., Vinyals, O., et al.: Neural discrete representation learning. *Advances in neural information processing systems* **30** (2017) [3](#), [4](#)
21. Wang, J., Zhang, M.: Deepflash: An efficient network for learning-based medical image registration. In: CVPR. pp. 4444–4452 (2020) [2](#)
22. Xu, J., Chen, E.Z., Chen, X., Chen, T., Sun, S.: Multi-scale neural odes for 3d medical image registration. In: MICCAI. pp. 213–223. Springer (2021) [2](#)
23. Xu, Z., Niethammer, M.: Deepatlas: Joint semi-supervised learning of image registration and segmentation. In: MICCAI. pp. 420–429. Springer (2019) [5](#)
24. Yang, Q., Fu, Y., Giganti, F., Ghavami, N., Chen, Q., Noble, J.A., Vercauteren, T., Barratt, D., Hu, Y.: Longitudinal image registration with temporal-order and subject-specificity discrimination. In: MICCAI. pp. 243–252. Springer (2020) [2](#), [7](#)
25. Ye, M., Kanski, M., Yang, D., Chang, Q., Yan, Z., Huang, Q., Axel, L., Metaxas, D.: Deeptag: An unsupervised deep learning method for motion tracking on cardiac tagging magnetic resonance images. In: CVPR. pp. 7261–7271 (June 2021) [8](#)
26. Zeng, Q., Fu, Y., Tian, Z., Lei, Y., Zhang, Y., Wang, T., Mao, H., Liu, T., Curran, W.J., Jani, A.B., et al.: Label-driven magnetic resonance imaging (mri)-transrectal ultrasound (trus) registration using weakly supervised learning for mri-guided prostate radiotherapy. *Physics in Medicine & Biology* **65**(13), 135002 (2020) [5](#)

27. Zhang, M., Liao, R., Dalca, A.V., Turk, E.A., Luo, J., Grant, P.E., Golland, P.: Frequency diffeomorphisms for efficient image registration. In: IPMI. pp. 559–570. Springer (2017) [2](#)
28. Zhang, Y., Pei, Y., Zha, H.: Learning dual transformer network for diffeomorphic registration. In: MICCAI. pp. 129–138. Springer (2021) [2](#)