001

002

003

004 005

006

007

008

009 010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

054 055 056 057 058 059 060 061 062 063 064 065 066 067 068 069 070 071 072 073 074 075 076 077 078 079 080 081 082 083 084 085 086 087 088 089 090 091 092 093 094 095 096 097 098 099 100 101 102 103 104 105 106 107

FOAL: Fast Online Adaptive Learning for Cardiac Motion Estimation

Anonymous CVPR submission

Paper ID 4635

Abstract

Motion estimation of cardiac MRI videos is crucial for the evaluation of human heart anatomy and function. Recent researches show promising results with deep learning based methods. In clinical deployment, however, they suffer dramatic performance drops due to mismatched distributions between training and testing datasets, commonly encountered in the clinical environment. On the other hand, it is arguably impossible to collect all representative datasets and to train a universal tracker before deployment. In this context, we proposed a novel fast online adaptive learning (FOAL) framework: an online gradient descent based optimizer that is optimized by a meta-learner. The metalearner enables the online optimizer to perform a fast and robust adaptation. We evaluated our method through extensive experiments on two public clinical datasets. The results showed the superior performance of FOAL in accuracy compared to the offline-trained tracking method. On average, the FOAL took only 0.4 second per video for online optimization.

1. Introduction

036 Video dense tracking and motion estimation using deep 037 learning has gained great progress for natural image appli-038 cations in recent research [34, 12, 21, 43, 45, 16, 11, 39, 18, 039 26, 47, 20]. In medical imaging, videos compared to static 040 images, are ideal for dynamically changing physiological 041 processes such as the beating heart and are commonly used in clinical settings. Feature tracking of dynamic cardiac im-042 043 ages can provide precise and comprehensive assessments of 044 the cardiac motion and has been proved valuable for car-045 diac disease management [33, 27, 41, 23]. Motion estimation can also benefit other tasks in cardiac imaging, such as 046 047 image reconstruction [10, 30] and semi-supervised segmen-048 tation [25, 37, 46, 17, 42]. Recently, deep learning based methods show promising results in cardiac motion estima-049 tion [25, 46, 15, 22]. However, most studies have been de-050 signed in a research environment: the proposed models are 051 052 trained and tested on the data with similar distributions. In 053 a clinical environment, however, the imaged objects may

present various anatomies (abnormally thin or thick heart muscle) and/or dynamics (irregularly beating heart) for different diseases. On top of that, the imaging process itself commonly introduces many, if not more, variations. This is especially true for cardiac magnetic resonance (CMR) imaging, which provides superior video quality over ultrasound, but the image appearances are influenced by multiple factors including scanner vendors, main magnetic fields, different scanning protocols and technicians' operations. It is arguably impossible to build a dataset that includes every combination of the variations and train a universal tracker on it. It is also not ideal and sometimes impossible in a clinical setting that the pre-trained network gets fine-tuned on the data from a different distribution, given the scarcity nature of medical data. In other words, for a clinically suitable deep-trained tracker, the neural network needs to possess the capability to quickly adapt to new data from unseen distributions. Towards this end, we propose a fast online adaptive learning (FOAL) mechanism for dense video tracking applied to cardiac motion estimation. The proposed framework consists of an online adaptive stage and an offline meta-learning stage. The offline meta-learning trains the model to gain the adaptation capability and the online stage will apply this adaptation to adjust the model parameters using very few and unseen data. We have designed a unique module for video tracking used in both stages to train an adaptive tracker. The tracker trained using the proposed FOAL achieves the state-of-the-art (SOTA) results compared to strong baselines. The contributions of our work are summarized as follows.

- In the context of dense motion estimation, we proposed a novel online model adaptation method, which adapts a trained baseline model to a new video using a gradient descent optimization.
- We proposed a meta-learning method optimizing the proposed online optimizer. The meta-learner enables the online optimizer to perform a fast and robust adaption.
- We proposed practical solutions for training meta learner in dense motion estimation task.

• Our proposed method is not limited to the network structure of the baseline dense motion estimation. The extensive experiments consistently demonstrated supe-

rior performance improvement of our method in accuracy comparing to the baseline model.



Figure 1. Typical cardiac cycle of a healthy subject. The cycle indicates heart relaxation and contraction process. The myocardium (MYO) appears as a dark ring in the image. The left ventricle (LV) is filled with hyperintense blood signal contained inside the ring. The right ventricle (RV) cavity is indicated via red line.

2. Related Work

Section 2.1 discusses state-of-the-arts in the literature for motion estimation in the computer vision field. Section 2.2 introduces the task of cardiac motion estimation and existing studies on this topic. Section 2.3 introduces the modelagnostic meta-learning which has inspired our method.

2.1. Motion Estimation for Camera Videos

Motion estimation is one of the fundamental problems in computer vision field. In the literature, there are a few deep learning based approaches solving motion estimation such as reported works in [2, 5, 12, 34, 21]. Dosovitskiy et al. [5] proposed two optical flow estimation networks (Flownets): FlownetSimple and FlownetCorr. The former is a generic architecture and the latter includes a correlation layer to fuse feature vectors at different image locations. Flownet 2.0 in the work [12] further adds an extra branch to deal with pairs with small displacement and uses the original Flownet to deal with large displacement. Sun et al. [34] proposed a smaller and more efficient neural network structure utilizing feature pyramid as well as cost volume to get a more accurate motion. Most of these above works used a supervised learning approach with true motion field. In contrast to these supervised methods, Meister *et al.* [21] proposed an unsupervised framework where the flow was predicted and used to warp the source image to the reference image. The model is optimized to minimize the difference between warped image and reference image. In addition, an occlusion-aware forward-backward consistency loss is used with the census transform to improve the track-ing results. Note that our baseline model utilized a similar self-supervision idea as [21].

2.2. Cardiac Motion Estimation

CVPR 2020 Submission #4635. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Cardiac motion estimation takes a time series (video) of CMR images as input and predicts the heart motion through time. Motion fields are usually estimated at a pixel level due to the non-rigid nature of cardiac contraction. Normally the video records a complete cardiac contraction cycle: from the onset of contraction (end diastole ED), then to maximum contraction (end systole ES) and back to relaxation. Fig. 1 shows example CMR frames from a video of a normal subject. Motion of a frame is usually estimated relative to a reference frame that is commonly chosen as the ED or ES frame. Let frame at time t be I(x, y, t), and $I(x, y, t_{ref})$ as the reference image. The goal of motion estimation is to find the mapping F_{θ} such that

$$F_{\theta}: (I(x, y, t_{ref}), I(x, y, t)) \longrightarrow V_x(x, y, t), V_y(x, y, t)$$
(1)

where F_{θ} is the mapping function with parameter θ and V_x, V_y are the motion fields along x and y directions, respectively. Motion tracking methods can be generally categorized according to the different formulations of F_{θ} : optical flow based, conventional image registration based, and deep learning based.

Optical flow based method is built on several presumptions on image appearance and motion strength, such as brightness consistency and small motion between source and reference frames. The problem of applying optical flow based methods to CMR motion estimation is that the presumptions are violated in CMR videos [6]. Fig. 2 shows example images, illustrating the challenges of CMR.



Figure 2. Examples of challenges in CMR motion estimation. (a) and (b) are from one CMR video, where the upper part of the LV myocardium (anterior wall) has a big intensity drop due to the changes in MR coil detection sensitivity. (c) and (d) are from another CMR video, where large motion occurs between an early systole frame (source) and an end systole frame (reference).

In addition to the optical flow based approaches, image registration based methods [24, 28, 4, 31, 32, 36, 15] were applied to solve cardiac motion estimation.Craene *et al.* [4] utilized a B-spline velocity fields with physical constraints to compute the trajectories of feature points and performed the tracking. Rueckert *et al.* [28] proposed a free form deformation (FFD) method solving a general deformable image registration problem and recent work [24, 31, 32, 36] utilize this method to estimate the cardiac motion. It is known that FFD-like methods suffer from computation effi-

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

3. Method

We proposed an online adaptive tracking framework in the context of the dense motion tracking utilizing a deep neural network. The proposed method is a general video tracking framework which is not limited to motion estimation in CMR. Nevertheless, without loss of generality the method is presented in the CMR context.

MAML method inspired us to utilize a meta learner which

teaches the model to learn how to adapt to a new video.

ciency problem. To address this issue, Vigneault *et al.* [38] proposed a coarse-to-fine registration framework to track cardiac boundary points. This solution improved the time efficiency but an extra segmentation step was required. In addition, this sparse tracking lost motion understanding in the heart muscle region.

Recent success in deep neural network solving many computer vision problems has inspired efforts to explore deep learning based cardiac motion estimation. Qin et al. [25] proposed a multi-task framework that combines segmentation and motion estimation tasks. The learned cardiac motion field is used to warp the segmentation mask and guide the segmentation module in a semi-supervised manner. The results show that both segmentation and motion estimation performance are improved comparing to single task. Zheng et al. [46] proposed the apparent flow net which is a modified U-net. The segmentation masks were used in apparent flow net in order to improve the motion estimation. In work [15], a conditional variational autoencoder (VAE) based method was presented to estimate the cardiac motion. The VAE encoder is used to map deformations to latent variables, which is regularized via Gassian distribution and decode to a deformation filed via VAE decoder. Note that it is generally hard to obtain true cardiac motion and thus above works were quantitatively evaluated using the segmentation masks. In this work, we also use this type of evaluation.

2.3. Model Agnostic Meta Learning

Meta learning, or leaning to learn, aims to build a universal meta model that could make fast adaptation to new tasks [29]. Model-agnostic meta learning (MAML) [7] is a general strategy that searches for a good model-agnostic initialization parameters that are trained through training tasks and can quickly adapt to new tasks. Given the initial model parameters θ , for every task T_i in the training set, the taskspecific parameters θ_i are independently updated within the task dataset using gradient descent with a differentiable loss function *L*:

$$\theta_i \leftarrow \theta - \alpha \nabla_{\theta} L(T_i; \theta).$$
 (2)

Then the original model parameters θ are updated over all the training tasks:

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{i} L(T_i; \theta_i).$$
 (3)

Through these meta-training processes, the optimal "ini-262 263 tialization" parameters are supposed to be sensitive to 264 new task adaptation within limited number of adaptation steps. MAML has been widely used in few-shot learn-265 ing [8, 35, 9], neural architecture search [19], graphical neu-266 ral network [9], compressed sensing [40] and transfer learn-267 ing [44]. Most applications using MAML are to solve high-268 269 level vision tasks such as classification and recognition. The Figure 3. Overview of the dense tracking framework. The encoder is a Siamese structure that takes source and reference images as input. The feature maps produced by Siamese encoder are concatenated and fed into the decoder.

Dense Motion Tracke

3.1. Dense Motion Tracking

Fig. 3 depicts the architecture of our dense tracking framework. The overall idea of the dense motion tracking is an end-to-end unsupervised learning approach that inspired from [21]. Annotating motion field for heart is an intractable task and unsupervised learning avoids the necessity of the ground truth. In our work we used a lightweight backbone of the network: the inputs are source image and reference image (e.g. two frames in the same video). The encoder is a Siamese [3] structure. The decoder is a series of convolution and transpose convolution operators used to decode the features and restore the output to original image size. The output is the predicted motion field. To perform unsupervised learning, the spatial transformer network [13] is utilized to deform/warp the source image to the reference image and image reconstruction loss L_{mse} is used to minimize the difference between the warped source image and the reference image. L_{mse} is the mean square error (MSE). In addition to L_{mse} , motion field smoothness L_{smooth} proposed in [25] is used to avoid abrupt motion change and a bidirectional (forward-backward) flow consistency loss L_{con} proposed in [21] is used. The total loss L_{total} is thus defined as follows:

$$L_{total} = L_{mse} + \alpha_s L_{smooth} + \beta_c L_{con}, \qquad (4)$$

where α_s and β_c are used to balance three losses.

3.2. Online Optimizer

The unsupervised dense tracking (Section 3.1) mitigates the need of ground truth motion fields. However, the dis271 272

270

273 274

275 276

277 278 279

280

281

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

343

344

345

346

347

361

362

363

364

365

366

367

368

369

370

371

372

373

374

398

399

400

401

402

403

404

405

406

407

408

409

410 411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

324 tribution mismatch between training and test datasets is a 325 continuous challenge, particularly the long tail problem in 326 medical image domain. The clinical deployment of a deep 327 learning model suffers domain mismatch problem. It is a 328 challenge to collect sufficient samples for training a uni-329 versal tracker. In this section, in the context of the pro-330 posed dense tracking, we extend the tracker to address the 331 dataset distribution mismatch problem. Instead of training 332 such a universal tracker offline, we make the tracker be-333 ing aware of the test data online. The idea behind this is 334 to enable a given tracker automatically adapting the model 335 to a new video x. Suppose we have a model f_{θ} using the 336 proposed dense tracker trained on dataset D_a with a distri-337 bution $p(D_a)$. The online adaptive learning on video x is 338 an online optimization algorithm and is realized via back-339 propagating through the stochastic gradient descent steps as 340 follows: 341

$$\theta' \leftarrow \theta' - \alpha \nabla_{\theta'} L(f_{\theta'}),$$
 (5)

where θ' represents the model parameters and is initialized from θ . α is the learning rate. We utilized the same loss function *L* defined in Eq.(4). The overview of the online adaptive algorithm is outlined in the Algorithm 1.

-	Algorithm 1 FOAL online optimization
1	
	Input: Single video t: x_t , learning rate: α , trained model:
	f_{θ} , number of online tracking optimization steps: m
	$ heta_t' \leftarrow heta$
	Sample K pairs $D_t = \{a_t^{(j)}, b_t^{(j)}\}$ from video x_t
	for t from 1 to m do
	Evaluate loss $L_t(f_{\theta'_t})$ using D_t
	Compute parameters with gradient descent:
	$\theta_t' \leftarrow \theta_t' - \alpha \nabla_{\theta_t'} L_t(f_{\theta_t'})$
	end for
(Output: updated network weights: θ'_t

It's worth pointing out that the gradient descent steps are performed over all parameters of the network at online stage. Thus, it is computationally expensive to optimize them on all image pairs (source and reference) with too many steps. We aim to adapt the offline model in just a few steps using only a small number of online samples. We realize this by employing meta-learning to optimize this optimization procedure. This idea is inspired by MAML [7], which is used to learn good initial model parameters via meta learning. Like in MAML, we perform a second order optimization by back-propagation using stochastic gradient descent through the online optimization Eq. (5).

3.3. Meta-learning

We utilized a meta leaner to re-train the model f_{θ} on the dataset D_{meta} from parameters θ in order to teach the online optimizer in Eq. (5) learn to adapt f_{θ} to a given video

Algorithm 2 FOAL offline meta-learning	378
Input: video set: X, learning rate: α , β , initial model: f_{θ} ,	37
number of online tracking optimization steps: m	381
while not done do	38
Sample N videos $\{x_1, x_2,, x_N\}$ from X	382
for i from 1 to N do	38:
$ heta_i' \leftarrow heta$	384
Sample K pairs $D_i = \int a^{(j)} b^{(j)} from video x$	385
Sample K pairs $D_i = \{u_i, b_i\}$ from video x_i	386
lor t from 1 to m do	387
Evaluate loss $L_i(f_{\theta'_i})$ using D_i	388
Compute parameters with gradient descent:	389
$\theta_i' \leftarrow \theta_i' - \alpha \nabla_{\theta_i'} L_i(f_{\theta_i'})$	39
end for	39
Sample K pairs $D'_i = \{a_i^{(k)}, b_i^{(k)}\}$ from video x_i	393
end for	30'
Model update: $\theta \leftarrow \theta - \beta \nabla_{\theta} \frac{1}{N} \sum_{i}^{N} L_{i}(f_{\theta'_{i}})$ using	20/
each D'_i and video-specific loss $L_i(f_{\theta'})$	394
end while	39:
Output: updated model θ	39
ourput apauta mouro	397

x. Note that D_{meta} is either $p(D_a)$ or a new distribution $p(D_b)$, where D_b is a new dataset, and $p(D_b)$ may mismatch domain $p(D_a)$. The full algorithm is outlined in Algorithm 2. There are two For-loops in Algorithm 2. The inner For-loop is the proposed optimization algorithm in Algorithm 1 for optimizing the online optimizer Eq. (5). The outer For-loop is the meta-leaner and the meta optimizer is defined as follows.

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \frac{1}{N} \sum_{i}^{N} L_{i}(f_{\theta_{i}'}),$$
(6)

where *i* is i^{th} video in the training procedure. *N* is the number of videos in a batch size for optimizing the meta learner. β is the learning rate of the meta-learner. L_i is the loss (Eq. 4) evaluating on the i^{th} video. $f_{\theta'_i}$ is the model parameters for the i^{th} video.

3.4. Practical Version of the Meta-Learning

Memory limitation and solution: In contrast to few-shot learning (a classification problem) discussed in MAML [7], dense motion tracker need store a larger number of feature maps (i.e. requiring a large amount of GPU memory) given a larger image size (e.g. 192×192). The meta optimizer (Eq. 6) requires computing derivatives of each independent model associated with a specific video. To tackle this problem, by employing the property that the gradient operator and the average operator are commutative in Eq. 6, we swap the two operators as shown in Eq. (7).

$$\nabla_{\theta} \frac{1}{N} \sum_{i}^{N} L_{i}(f_{\theta_{i}'}) \Leftrightarrow \frac{1}{N} \sum_{i}^{N} \nabla_{\theta} L_{i}(f_{\theta_{i}'}) \tag{7}$$

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

which enables computing gradients on GPU and transferring them to CPU.

First order derivative approximation: Note that in Eq. (7), second-order derivative is needed in backpropagation. This involves calculating second-order Hessian matrix, which is computationally costly. As a workaround, we use first-order approximation, whose effectiveness is demonstrated in MAML [7]. In [7], the approximation rendered comparable results to the second-order derivatives.

Table 1. Inside distribution v.s. outside distribution Dice coefficient results for baseline model, proposed FOAL without meta learning (FOAL w/o meta) and proposed FOAL with metalearning (FOAL + meta). Averaged Dice coefficient with standard deviation are given among five-fold leave-one-disease-out cross validation.

Method	LV	RV	MYO	
method	Inside Distribution Test Set			
Baseline	0.838(0.024)	0.825(0.013)	0.797(0.014)	
FOAL w/o meta	0.856(0.021)	0.842(0.013)	0.820(0.008)	
FOAL + meta	0.873(0.019)	0.859(0.013)	0.840(0.007)	
	Outside Distribution Test Set			
Baseline	0.840(0.094)	0.775(0.096)	0.803(0.045)	
FOAL w/o meta	0.863(0.077)	0.801(0.085)	0.828(0.031)	
FOAL + meta	0.880(0.065)	0.806(0.086)	0.846(0.027)	

Table 2. Inside distribution v.s. outside distribution Hausdorff distance (mm) results for baseline model, proposed FOAL without meta learning (FOAL w/o meta) and proposed FOAL with metalearning (FOAL + meta). Averaged Hausdorff distance with standard deviation are given among five-fold leave-one-disease-out cross validation.

Method	LV	RV	MYO		
	Inside Distribution Test Set				
Baseline	7.265(0.779)	8.782(0.422)	6.930(0.548)		
FOAL w/o meta	6.417(0.627)	8.141(0.329)	6.286(0.469)		
FOAL + meta	6.012(0.580)	7.731(0.303)	6.157(0.489)		
	Outside Distribution Test Set				
Baseline	6.921(2.147)	10.173(1.436)	6.716(1.803)		
FOAL w/o meta	6.158(1.727)	9.320(1.422)	6.107(1.506)		
FOAL + meta	5.832(1.534)	9.378(1.417)	5.987(1.437)		

4. Evaluation Methodology

In this section, we present evaluation methodology on compared tracking methods: tracking performed using proposed dense motion tracking method (baseline model), tracking performed using online optimization from the baseline model without meta-learning (FOAL without meta-learning), and tracking performed using online optimization with meta-learning (FOAL with meta-learning).

4.1. Datasets and Evaluation Reference

In our study, two public CMR datasets were utilized: ACDC dataset [1] and Kaggle Data Science Bowl Cardiac Challenge Data [14]. All data acquisitions were performed using breathholding so that only cardiac motion is observed in the videos. It is arguably impossible to make independent reference standard of the cardiac motion manually. To perform quantitative analysis, we utilized segmentation masks as the independent reference standard. In the test dataset of the study, we have heart segmentation references at both the first frame and the evaluated reference frame. We generate the segmentation masks via warping source segmentation to the reference and compare it to the annotation using quantitative indices defined in section 4.4.

ACDC Dataset: It includes short-axis view CMR videos from 100 subjects (healthy and diseased cases). Each subject contains multiple slices (9-10) and each slice is a video sequence covering at least one heart beat cycle. Overall, there are 951 videos in this dataset. Each video provides two heart segmentation masks: one for ED phase and one for ES phase. The segmentation labels are: right ventricle (RV) cavity, myocardium (MYO) and left ventricle (LV) cavity. In addition, the 100 subjects are evenly divided into 5 categories with 20 subjects each. These are diagnosed into: normal cases (NOR), systolic heart failure with infarction (MINF), dilated cardiomyopathy (DCM), hypertrophic cardiomyopathy (HCM), abnormal right ventricle (ARV). The CMR videos were collected over a 6 year period using two MRI scanners of different main magnetic fields: 1.5 T Siemens Area and 3.0 T Siemens Trio Tim (Siemens Medical Solutions, Germany) [1].

Kaggle Data Science Bowl Cardiac Challenge Dataset: It includes short-axis view CMR videos from 1100 subjects. Each subject contains multiple slices (8-10) and each slice is a video sequence covering at least one cardiac cycle. Overall, there are 11202 videos in this dataset. The original challenge is to predict ejection fraction from the videos. Ejection fraction ground truth was provided but irrelevant to our study. The subjects have a large health and age range and the images were collected from numerous sites [14]. However detailed information such as disease types is not disclosed nor there are segmentation labels. Nevertheless this large real clinical dataset can be used to train baseline dense motion model.

4.2. Implementation Details

For image preprocessing, we normalized the gray value to 0-255 and we applied center cropping and zero padding to adjust image size to 192×192 . All models are trained and tested on a Tesla V100 workstation. The other implementation details are presented as following.

Dense motion tracker: As for the baseline model, we

486 487

488

489

490

491

492

493

494

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

CVPR 2020 Submission #4635. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 4. The bar-plots of inside distribution v.s. outside distribution Dice coefficient results for baseline model, proposed FOAL without meta learning (FOAL w/o meta) and FOAL with meta-learning (FOAL with meta) for all five folds. Different diseases as outside distributions are presented in different columns. Top row is the inside distribution test and the bottom row is the outside distribution test. The data of the outside distribution disease were excluded in the baseline training and meta-training. Averaged values and standard deviations are presented.

adopted a lightweight (shallower and narrower) version of the motion prediction network proposed by Qin *et al.* [25]. We halved the number of feature maps of each layer and the number of layers. We set $\alpha_s = 5 \times 10^{-5}$ and $\beta_c = 10^{-6}$ in Eq. (4). The batch size is 20 images. We utilized Adam optimizer with an initial learning rate 10^{-3} .

Online optimizer: The number of update steps m = 3 and 571 the number of sampled pairs K = 24 in Algorithm 1. We 572 used Adam optimizer with learning rate $\alpha = 10^{-4}$.

Meta learner: We used the number of sampled videos 574 n = 2, the number of update steps m = 5, and the number 575 of sampled pairs K = 24 in the online optimization in Al-576 gorithm 2. SGD optimizer is used for online optimizer with 577 a fixed learning rate $\alpha = 10^{-5}$. Adam optimizer is used for 578 the meta-learner with an initial learning rate $\beta = 10^{-5}$ in 579 Algorithm 2. The meta training steps are 6,000.

4.3. Experiment Setups

Inside distribution vs Outside distribution: In data driven machine learning, we always hypothesize that train-ing samples and testing samples are drawn from the the same distribution (inside distribution). The violation of the hypothesis (outside distribution in the testing set) usually gives poor model generalization on testing set. In this study, we performed five fold cross validations in light of leaving-one-disease-out method on the ACDC dataset. The idea be-hind this is to separate inside distribution (P_{in}) and out-side distribution (P_{out}) in terms of known diseases. Due to the significant cardiac anatomy and dynamic differences

between different diseases, one disease category could be viewed as an outside distribution compared to the other 4 diseases. For subjects in the inside distribution set, we separate them into train set ($80 \times 80\% = 64$ subjects) as $p(D_a)$ and $p(D_{meta})$, and test set ($80 \times 20\% = 16$ subjects) as $p(D_{tinside})$. 100% subjects in the outside distribution (20 subjects) set were used in the test set as $p(D_{toutside})$. In this experiment, we trained and evaluated all three compared methods on the ACDC dataset.

Fine-tuning and Generalization: We observed that the proposed FOAL with meta-learning needs to train the metalearner from a baseline model. In the dense tracking context, it is difficult to train the meta-learner from scratch. However, our idea behind the FOAL is to enable any dense tracker to boost their performance via online optimization through meta-learning. To validate the generalizability, we utilized Kaggle dataset that is without any meta information. Specifically, we used the 30% subjects of the entire Kaggle dataset as $p(D_a)$ to train the baseline model. We then performed leave-one-disease-out cross-validation on the ACDC dataset. Note that the Kaggle data are only used for training the baseline model while $p(D_{meta})$, $p(D_{t_{inside}}) \mbox{ and } p(D_{t_{outside}}) \mbox{ are all from ACDC with the}$ same split in the first experiment. In addition to the leaveone-disease-out cross validation, starting from the baseline model trained on Kaggle, we also evaluated a vanilla finetuning model to FOAL with meta learning using 20% of the entire ACDC dataset $(100 \times 20\% = 20 \text{ subjects})$. 100%or 10% of the rest ACDC data were used to train the two models. All 5 categories were mixed.

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

Note that the vanilla fine-tuning model used the same training parameters as the baseline model except that we changed the learning rate to 10^{-5} .

4.4. Quantitative Metrics

We used the DICE coefficient (Eq. (8)) and Hausdorff distance error (Eq. (9)) as quantitative metrics to evaluate the compared tracking methods on segmentation masks. The metrics are defined as:

$$DICE = \frac{2 \times |S_A \cap S_B|}{|S_A| + |S_B|},\tag{8}$$

where S_A and S_B are the segmentation mask A and the segmentation mask B, respectively.

$$H(C_A, C_B) = \max_{a \in C_A} \{ \min_{b \in C_B} ||a - b||_2 \},$$
(9)

where *a* and *b* are the points on the contour A and the contour B, respectively. $|| \cdot ||_2$ is the Euclidean distance.

5. Results and Discussion

Inside distribution vs outside distribution on ACDC data: The five-fold cross-validation experiment in this part is described in Section 4.3. Fig. 4 depicts all three compared methods (baseline model, FOAL without meta learning and FOAL with meta learning) in every cross-validation with test samples drawn from inside or outside distribution. Table 1 and Table 2 summarize Dice and Hausdorff distance results, respectively, for both inside and outside distributions averaged over the five folds. Fig. 4, Table 1 and Table 2 show that the proposed FOAL with meta learning approach outperforms the baseline tracker. For the inside distribution test, our FOAL with meta learning increased the Dice by 3.7% and reduced Hausdorff distance error by 1.0 mm on average. It is worth pointing out that even the training and testing are within the same disease distribution, the variations from patients, scanner types, scanner settings etc. are still large, which can explain the reduced errors from our method compared to the baseline. The largest accuracy improvement occurs on MYO with 4.3% on Dice for both inside distribution and outside distribution. On the zero-shot (outside distribution) dataset, our FOAL with meta learning achieves superior performance (e.g. on average 3.8% increase on Dice) compared to the baseline. In 692 addition, we observed that FOAL with meta learning out-693 performs FOAL without meta learning consistently. This 694 695 demonstrates the effectiveness of the meta learning to en-696 hance the adaptation capability of the online optimizer. This result is not surprising because the online optimizer learns 697 how to adapt to a new video using offline meta training on 698 a large number of videos. This capability teaches the online 699 700 optimizer to find a sub-optimal path to a better solution than 701 the optimizer without meta learning can.

Fig. 5 depicts the warped segmentation results using corresponding deformation fields which were generated by baseline model and FOAL with meta-learning. In Fig. 5, ED and ES frames in the video are also illustrated. We observed a significant appearance and shape difference inside the heart region. Refering to annotations, our method improved LV (blue color) and MYO (green color) comparing to the baseline method.

Table 3. Finetuning experiment with Kaggle baseline training and ACDC inside and outside distribution test sets. Dice coefficients are averaged over the five-fold cross validation for baseline model trained on Kaggle data (Baseline), fine-tunined model on the ACDC dataset and FOAL with meta-learning (FOAL + meta) on the ACDC dataset.

Method	LV	RV	MYO	
	Inside Distribution Test Set			
Baseline	0.864(0.019)	0.847(0.013)	0.830(0.010)	
Finetune	0.861(0.023)	0.850(0.012)	0.827(0.014)	
FOAL + meta	0.880(0.017)	0.866(0.010)	0.847(0.009)	
	Outside Distribution Test Set			
Baseline	0.874(0.070)	0.796(0.093)	0.841(0.024)	
Finetune	0.870(0.070)	0.792(0.094)	0.833(0.031)	
FOAL + meta	0.885(0.059)	0.804(0.091)	0.849(0.023)	

Fine-tuning and Generalization: The experiment setup in this part is discussed in Section 4.3. We compared the baseline model trained on Kaggle data (Baseline), a model fine-tuned on ACDC data from the baseline model (Finetune) and our proposed FOAL with meta-learning from the baseline model (FOAL+meta). Averaged Dice coefficients among five folds for both inside distribution and outside distribution can be found in Table 3. The baseline model performs comparably well on both distributions except RV. This might be because Kaggle dataset consists of a variety of cardiac diseases and it has distribution overlaps with both the inside distribution and the outside distribution datasets but not for RV. Fine-tuning the model on the ACDC dataset does not improve the performance. Comparing to the baseline model, our method improved 2.7% on the inside distribution test and 2.4% on the outside distribution test in terms of Dice.

Table. 4 shows Dice results for vanilla fine-tuning model and our FOAL with meta-learning using 10% or 100% ACDC training samples. In contrast to the leave-onedisease-out experiments, we did not isolate any disease in the training samples in this experiment and the models were tested on the entire ACDC test set. Vanilla fine-tuning model made the performance slightly worse in the 10% experiment while it slightly improved the accuracy in the 100% experiment comparing to the baseline model. Meanwhile, FOAL with meta-learning gave 1.68% and 1.71% Dice increases on average for both 10% and 100% exper-



Figure 5. Examples of the tracking results of the mask overlays warped from ED heart phase to ES heart phase. The warp operation utilized deformation fields which were generated from the compared methods. From left to right: the starting frame (ED phase), the final frame (ES phase), baseline model, FOAL with meta-learning and the expert mask annotations. Note that the red mask represents RV, green represents MYO and blue represents LV.

iments, respectively. This result is consistent with above five fold cross-validation test. In addition, Fig. 4 demonstrates that our FOAL performs comparably well using a small amount of data when it meta-trained from a strong baseline model.

On average, our FOAL online optimization algorithm requires 413 ± 8 milliseconds (mean \pm standard deviation), which we find it completely durable for most current clinical applications.

Table 4. Finetuning experiment with Kaggle basline training and 100% and 10% ACDC training dataset. Dice coefficients for baseline model trained on Kaggle data (Baseline), vanilla fine-tuned model on the ACDC (Finetune) and FOAL with meta-learning on the ACDC (FOAL + meta) are reported.

the nebe (10h	E + lifeta) are le	poneu.		
Method	LV	RV	MYO	
	100% of ACDC training data			
Baseline	0.865(0.103)	0.845(0.080)	0.829(0.065)	
Finetune	0.865(0.104)	0.854(0.079)	0.831(0.063)	
FOAL	0.881(0.086)	0.865(0.070)	0.845(0.051)	
	10% (of ACDC trainin	g data	
Baseline	0.865(0.103)	0.845(0.080)	0.829(0.065)	
Finetune	0.864(0.104)	0.845(0.082)	0.824(0.073)	
FOAL +meta	0.882(0.086)	0.863(0.071)	0.845(0.051)	

6. Conclusion

In this work, we proposed a novel online adaptive learning method to minimize the domain mismatch problem in the context of dense cardiac motion estimation. The online adaptor is a gradient descent based optimizer which itself is also optimized by a meta-learner. The meta-learning strategy allows the online optimizer to perform a fast adaption using a limited number of model updates and a small number of image pairs from a single video. The tracking performance is significantly improved in all the zero-shot (outside distribution comparing to the training samples) experimental setups. In addition, it is observed that the online adaptor is able to minimize the tracking errors in the inside distribution tests. Experimental results demonstrate that our methods obtain superior performance compared to the model without online adaption. The pilot study shows the feasibility of applying tin the context of unsupervised dense motion tracking or deformable image registration. The proposed method provides a practical and elegant approach to an often overlooked problem in existing art. We hope to inspire more discussions and work to benefit other clinical applications suffer from similar issues.

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

864 References

- [1] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: Is the problem solved? *IEEE transactions on medical imaging*, 37(11):2514–2525, 2018. 5
- [2] Jingchun Cheng, Yi-Hsuan Tsai, Shengjin Wang, and Ming-Hsuan Yang. Segflow: Joint learning for video object segmentation and optical flow. In *Proceedings of the IEEE international conference on computer vision*, pages 686–695, 2017. 2
 - [3] Sumit Chopra, Raia Hadsell, Yann LeCun, et al. Learning a similarity metric discriminatively, with application to face verification. In *CVPR* (1), pages 539–546, 2005. 3
- [4] Mathieu De Craene, Gemma Piella, Oscar Camara, Nicolas Duchateau, Etelvino Silva, Adelina Doltra, Jan D'hooge, Josep Brugada, Marta Sitges, and Alejandro F Frangi. Temporal diffeomorphic free-form deformation: Application to motion and strain estimation from 3d echocardiography. *Medical image analysis*, 16(2):427–450, 2012. 2
- [5] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer* vision, pages 2758–2766, 2015. 2
- [6] James P Earls, Vincent B Ho, Thomas K Foo, Ernesto Castillo, and Scott D Flamm. Cardiac mri: recent progress and continued challenges. *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 16(2):111–127, 2002. 2
 - [7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Modelagnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1126–1135. JMLR. org, 2017. 3, 4, 5
- [8] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via metalearning. arXiv preprint arXiv:1709.04905, 2017. 3
- [9] Victor Garcia and Joan Bruna. Few-shot learning with graph neural networks. *arXiv preprint arXiv:1711.04043*, 2017. 3
- [10] Qiaoying Huang, Dong Yang, Hui Qu, Jingru Yi, Pengxiang Wu, and Dimitris N Metaxas. Dynamic mri reconstruction with motion-guided network. 2018. 1
- [11] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), June 2019. 1
- [12] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2462–2470, 2017. 1, 2

- [13] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In Advances in neural information processing systems, pages 2017–2025, 2015. 3
- [14] kaggle. Data science bowl cardiac challenge data, 2014. Second Annual Data Science Bowl from kaggle, https://www.kaggle.com/c/ second-annual-data-science-bowl/data. 5
- [15] Julian Krebs, Hervé e Delingette, Boris Mailhé, Nicholas Ayache, and Tommaso Mansi. Learning a probabilistic model for diffeomorphic registration. *IEEE transactions on medical imaging*, 2019. 1, 2, 3
- [16] Hsueh-Ying Lai, Yi-Hsuan Tsai, and Wei-Chen Chiu. Bridging stereo matching and optical flow via spatiotemporal correspondence. In *The IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), June 2019. 1
- [17] Matthew Chung Hai Lee, Kersten Petersen, Nick Pawlowski, Ben Glocker, and Michiel Schaap. Tetris: Template transformer networks for image segmentation with shape priors. *IEEE transactions on medical imaging*, 2019. 1
- [18] Xin Lei, Liangyu He, Yixuan Tan, Ken Xingze Wang, Xinggang Wang, Yihan Du, Shanhui Fan, and Zongfu Yu. Direct object recognition without line-of-sight using optical coherence. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1
- [19] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. arXiv preprint arXiv:1806.09055, 2018. 3
- [20] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Selflow: Self-supervised learning of optical flow. In *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), June 2019. 1
- [21] Simon Meister, Junhwa Hur, and Stefan Roth. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 1, 2, 3
- [22] Manuel A Morales, David Izquierdo-Garcia, Iman Aganj, Jayashree Kalpathy-Cramer, Bruce R Rosen, and Ciprian Catana. Implementation and validation of a threedimensional cardiac motion estimation network. *Radiology: Artificial Intelligence*, 1(4):e180080, 2019. 1
- [23] Asif Padiyath, Paul Gribben, Joseph R Abraham, Ling Li, Sheela Rangamani, Andreas Schuster, David A Danford, Gianni Pedrizzetti, and Shelby Kutty. Echocardiography and cardiac magnetic resonance-based feature tracking in the assessment of myocardial mechanics in tetralogy of fallot: an intermodality comparison. *Echocardiography*, 30(2):203– 210, 2013. 1
- [24] Esther Puyol-Antón, Bram Ruijsink, Wenjia Bai, Hélène Langet, Mathieu De Craene, Julia A Schnabel, Paolo Piro, Andrew P King, and Matthew Sinclair. Fully automated myocardial strain estimation from cine mri using convolutional neural networks. In 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018), pages 1139–1143. IEEE, 2018. 2
- [25] Chen Qin, Wenjia Bai, Jo Schlemper, Steffen E Petersen, Stefan K Piechnik, Stefan Neubauer, and Daniel Rueckert.

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

972Joint learning of motion estimation and segmentation for car-973diac mr image sequences. In International Conference on974Medical Image Computing and Computer-Assisted Interven-975tion, pages 472–480. Springer, 2018. 1, 3, 6

- [26] Anurag Ranjan, Varun Jampani, Lukas Balles, Kihwan Kim,
 Deqing Sun, Jonas Wulff, and Michael J. Black. Competitive
 collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1
- [27] Martin Reindl, Christina Tiller, Magdalena Holzknecht, Ivan Lechner, Alexander Beck, David Plappert, Michelle Gorzala, Mathias Pamminger, Agnes Mayr, Gert Klug, et al. Prognostic implications of global longitudinal strain by featuretracking cardiac magnetic resonance in st-elevation myocardial infarction. *Circulation: Cardiovascular Imaging*, 12(11):e009404, 2019. 1
- [28] Daniel Rueckert, Luke I Sonoda, Carmel Hayes, Derek LG
 Hill, Martin O Leach, and David J Hawkes. Nonrigid
 registration using free-form deformations: application to
 breast mr images. *IEEE transactions on medical imaging*,
 18(8):712–721, 1999. 2
- [30] Gavin Seegoolam, Jo Schlemper, Chen Qin, Anthony Price, Jo Hajnal, and Daniel Rueckert. Exploiting motion for deep learning reconstruction of extremely-undersampled dynamic mri. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 704–712.
 Springer, 2019. 1
- [31] Dinggang Shen, Hari Sundar, Zhong Xue, Yong Fan, and Harold Litt. Consistent estimation of cardiac motions by 4d image registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 902–910. Springer, 2005. 2
- [32] Wenzhe Shi, Xiahai Zhuang, Haiyan Wang, Simon Duckett, Duy VN Luong, Catalina Tobon-Gomez, KaiPin Tung, Philip J Edwards, Kawal S Rhode, Reza S Razavi, et al. A comprehensive cardiac motion estimation framework using both untagged and 3-d tagged mr images based on non-rigid registration. *IEEE transactions on medical imaging*, 31(6):1263–1275, 2012. 2
- [33] Nicholas B Spath, Miquel Gomez, Russell J Everett, Scott
 Semple, Calvin WL Chin, Audrey C White, Alan G Japp,
 David E Newby, and Marc R Dweck. Global longitudinal
 strain analysis using cardiac mri in aortic stenosis: Comparison with left ventricular remodeling, myocardial fibrosis, and
 2-year clinical outcomes. *Radiology: Cardiothoracic Imag- ing*, 1(4):e190027, 2019. 1
- [34] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8934– 8943, 2018. 1, 2
- [35] Qianru Sun, Xinzhe Li, Yaoyao Liu, Shibao Zheng, Tat Seng Chua, and Bernt Schiele. Learning to self-train

for semi-supervised few-shot classification. *arXiv preprint arXiv:1906.00562*, 2019. **3**

- [36] Catalina Tobon-Gomez, Mathieu De Craene, Kristin Mcleod, Lennart Tautz, Wenzhe Shi, Anja Hennemuth, Adityo Prakosa, Hengui Wang, Gerry Carr-White, Stam Kapetanakis, et al. Benchmarking framework for myocardial tracking and deformation algorithms: An open access database. *Medical image analysis*, 17(6):632–648, 2013. 2
- [37] Gabriele Valvano, Agisilaos Chartsias, Andrea Leo, and Sotirios A Tsaftaris. Temporal consistency objectives regularize the learning of disentangled representations. In *Domain Adaptation and Representation Transfer and Medical Image Learning with Less Labels and Imperfect Data*, pages 11–19. Springer, 2019. 1
- [38] Davis M Vigneault, Weidi Xie, David A Bluemke, and J Alison Noble. Feature tracking cardiac magnetic resonance via deep learning and spline optimization. In *International Conference on Functional Imaging and Modeling of the Heart*, pages 183–194. Springer, 2017. 3
- [39] Yang Wang, Peng Wang, Zhenheng Yang, Chenxu Luo, Yi Yang, and Wei Xu. Unos: Unified unsupervised opticalflow and stereo-depth estimation by watching videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1
- [40] Yan Wu, Mihaela Rosca, and Timothy Lillicrap. Deep compressed sensing. arXiv preprint arXiv:1905.06723, 2019. 3
- [41] R Xia, T Zhu, Y Zhang, YS Chen, L Wang, JC Liao, YM Li, FJ Lü, and FB Gao. Tracking early reperfused myocardial infarction using cardiac mr. *Sichuan da xue xue bao*. *Yi xue ban= Journal of Sichuan University. Medical science edition*, 50(4):489–493, 2019. 1
- [42] Fan Yang, Yan Zhang, Pinggui Lei, Lihui Wang, Yuehong Miao, Hong Xie, and Zhu Zeng. A deep learning segmentation approach in free-breathing real-time cardiac magnetic resonance imaging. *BioMed research international*, 2019, 2019. 1
- [43] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1983–1992, 2018. 1
- [44] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3712–3722, 2018. 3
- [45] Mingliang Zhai, Xuezhi Xiang, Rongfang Zhang, Ning Lv, and Abdulmotaleb El Saddik. Optical flow estimation using dual self-attention pyramid networks. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019. 1
- [46] Qiao Zheng, Hervé Delingette, and Nicholas Ayache. Explainable cardiac pathology classification on cine mri with motion characterization by semi-supervised learning of apparent flow. *Medical image analysis*, 2019. **1**, **3**
- [47] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based learning of optical flow, depth, and egomotion. In *The IEEE Conference* on Computer Vision and Pattern Recognition (CVPR), June 2019. 1