

Towards Predicting Fine Finger Motions from Ultrasound Images via Kinematic Representation

Dean Zadok¹, Oren Salzman¹, Alon Wolf² and Alex M. Bronstein¹

Abstract—A central challenge in building robotic prostheses is the creation of a sensor-based system able to read physiological signals from the lower limb and instruct a robotic hand to perform various tasks. Existing systems typically perform discrete gestures such as pointing or grasping, by employing electromyography (EMG) or ultrasound (US) technologies to analyze muscle states. While estimating finger gestures has been done in the past by detecting prominent gestures, we are interested in detection, or inference, done in the context of fine motions that evolve over time. Examples include motions occurring when performing fine and dexterous tasks such as keyboard typing or piano playing. We consider this task as an important step towards higher adoption rates of robotic prostheses among arm amputees, as it has the potential to dramatically increase functionality in performing daily tasks. To this end, we present an end-to-end robotic system, which can successfully infer fine finger motions. This is achieved by modeling the hand as a robotic manipulator and using it as an intermediate representation to encode muscles' dynamics from a sequence of US images. We evaluated our method by collecting data from a group of subjects and demonstrating how it can be used to replay music played or text typed. To the best of our knowledge, this is the first study demonstrating these downstream tasks within an end-to-end system.

I. INTRODUCTION

Bionic hands or robotic prostheses were sought after for decades yet the first commercially-available robotic prostheses only became available in 2007 [1] and a fully-functioning prosthesis enabling all range of daily tasks is still out of reach. Existing robotic prostheses are typically connected to the limb of an amputee (or to a subject with severe phalanx or palm deformations) when the subject's muscles still respond to electric potentials generated by the brain [2]. These electric potentials allow inferring different finger motions by placing a non-invasive sensor on the residual limb, having subjects attempt these motions or gestures, and then discriminating between the different types of recorded electrical signals [3], [4]. This sensing technology, called electromyography (EMG), allows the reproduction of grasp flexion and extension yet is unable to reproduce finer motions typically involved in daily tasks.

Recently, it was shown that replacing EMG with ultrasound (US) imaging that captures the muscles' morphological state allows for better differentiation between discrete gestures or to classify full-finger flexion [5]. These images are obtained by placing an US probe on the residual limb and its efficacy is based on the fact that the muscles generate

different deformation patterns for different actions, and this behavior can be generalized among different subjects [5], [6], [7], [8]. Arguably, the most popular approach to infer motions from (high-dimensional) US images is via data-driven methods [6], [8], [9]. However, to the best of our knowledge, these methods still fall short of fully predicting dexterous hand motions or fine finger motions that a subject intended to perform. Predicting such motions, which is the focus of our work, is instrumental to the fine control of robotic hands and, consequently, to increase adoption rates of robotic prostheses.

Specifically, the key question this research aims to answer is “*To what extent can we predict and differentiate between fine finger motions by only having access to the lower-arm muscles?*” As challenging exemplary use cases, we consider piano playing and keyboard typing. In these tasks, finger gestures required to press or type are usually finer than those required to grasp objects or fully flex the fingers, generating smaller lower-arm muscle movements. To the best of our knowledge, predicting finger motions for these tasks has not been evaluated as part of an end-to-end system and we hope that our method and dataset will be used as a baseline for understanding these motions.

To this end, we present an end-to-end robotic system that allows to successfully solve our research question. This success stems from the design choices we took to construct our system which we view as a major contribution of this work. These include: (i) choosing an US sensor to capture the muscle's morphological state, (ii) using a sequence of multiple US images at each step, and (iii) a learning framework designed to exploit the temporal nature of the input US images. The proposed neural network (NN) architecture allows to inject domain knowledge accounting for the fact that pressing can be inferred by both the muscles' dynamics and an intermediate lower-dimensional representation of the palm of the hand inspired by robotic manipulators.

As we demonstrate in our empirical evaluation our framework can be used to accurately reproduce music played or text typed, applications that were previously considered unattainable. This was done by collecting data from a group of subjects and successfully predicting finger motions in our two motivating applications.¹ Finally, we demonstrate our system by connecting it to a robotic arm that replayed piano notes played by a subject in real-time.

¹Department of Computer Science, Technion, Haifa, Israel {deanzadok,osalzman,bron}@cs.technion.ac.il

²Department of Mechanical Engineering, Technion, Haifa, Israel alonw@me.technion.ac.il

¹Note that we do not showcase our results on deformed muscles. We performed tests on healthy subjects with fully-functioning hands in order to provide an accurate benchmark for the task we concentrate on. Our code and dataset are available on <https://github.com/deanzadok/finemotions>.

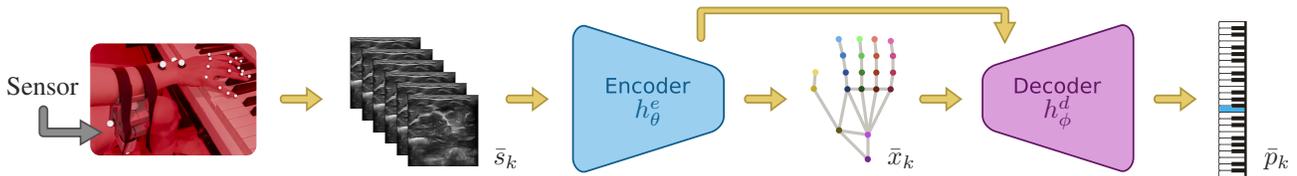


Fig. 1: A schematic flow of the proposed model-based method (Sec. V-C). An ultrasound sensor is placed on the lower arm while the subject is playing the piano. A continuous stream of ultrasound images \bar{s}_ℓ from the sensor is fed into the neural-network encoder h_θ^e which creates a latent representation of hand skeleton configurations \bar{x}_ℓ . The decoder h_ϕ^d receives the latent representation and outputs the vector of probabilities \bar{p}_ℓ indicating the pressed keys. The entire system is trained to produce a common representation for both the skeleton tracking and key prediction tasks.

II. RELATED WORK

Object grasping is one of many actions human hands perform daily. One reason for the low adoption rates of prosthetic arms among arm amputees is the lack of functionality in performing daily tasks [10]. To enrich the variance of gestures, advanced wearable surface-EMG (sEMG) sensors have been used (e.g., Myo armband [4], [11]). The array of features provided by sEMG enabled inferring more than grasp intention, and classifying a varied set of hand gestures via machine learning (ML) and deep learning (DL) [11], [12]. Additionally, sEMG demonstrated promising performance for arm amputees [4]. Noteworthy, in this domain, subjects differ in their gestures, which puts forth the requirement of operating across a variety of muscular structures [13].

In recent years, US sensing technologies emerged as a promising alternative to EMG. McIntosh et al. proposed to extract the optical flow of an US stream prior to the classification [9]. Others used ML algorithms demonstrating gesture classification from downsampled US images of the same area [5], [14]. Gesture recognition was achieved with wearable single-transducer sensors [6], [15]. In addition, researchers demonstrated that similar techniques allow decoupling of different degrees of freedom (DOFs) and use them to classify gestures [8], [16]. Unsupervised techniques were also studied, showing discrimination of unlabeled gestures [7], which might imply the possibility to understand muscle behavior of arm amputees without explicit labeling. Nevertheless, the feasibility of instructing a robotic hand to execute fine finger motions remains an open question.

For the analysis of US imagery DL techniques have become increasingly dominant in tasks ranging from tissue segmentation and pathology detection [17], and image enhancement [18], replacing traditional beamforming techniques with learned beamformers for faster and better image acquisition [19], [20]. In the context of learning dynamics of human organs, US imaging with transducers placed under the jaw demonstrated reliable speech and pronunciation recreation from imaged tongue movements [21], [22]. Vogt et al. utilized the tongue’s morphological behavior to extrapolate sung pitch and phonemes [23]. Lulich et al. studied the correlation between tongue motion and the frequencies generated by clarinet playing [24]. These works suggest that US provides a rich signal source describing minute muscular motion.

III. ANATOMICAL BACKGROUND

In this section, we provide the anatomical background required to understand the approach we take to solve our inference problem. Specifically, we explain (in general terms) the relation between lower-arm muscles and the state of the wrist’s and fingers’ joints. The lower arm contains two main bones, the *Ulna* and the *Radius*. Surrounding these bones are different muscles that are responsible for the flexion, extension, and rotation of different joints in our palms (including all the fingers). Two annotated US images of this region are presented in Fig. 2a and Fig. 2b, highlighting the muscles relevant to our problem as well as different muscle states corresponding to different finger motions.

We start by concentrating on the four fingers (Index through Little finger). There are three types of joints associated with these four fingers: the metacarpophalangeal (MP) joints, the proximal interphalangeal (PIP) joints, and the distal interphalangeal (DIP) joints (Fig. 2c). Close to the skin shell, we can find the *flexor carpi radialis* (FCR) and the *flexor carpi ulnaris* (FCU) that are (together with smaller, less prominent muscles) in charge of flexion of the wrist. Between the FCU and the FCR muscles, we can identify the *flexor digitorum superficialis* (FDS) that flexes the four fingers, and, in full flexion, also contributes to the flexion of the wrist. The FDS connects to tendons that go through the PIP and MP joints of the four fingers. Deeper into the arm, the *flexor digitorum profundus* (FDP), similar to the FDS, is in charge of finger flexion and wrist-control assistance using tendons connected to the DIP joints. Lastly, close to the FDP, is the *flexor pollicis longus* (FPL), which controls the thumb joints. To summarize, observing the state of the FDS and FDP muscles allows (in theory) to estimate the whereabouts of all four fingers, while the FPL provides the knowledge required to understand thumb movement.

IV. SYSTEM DESIGN

Following the anatomical background presented, we argue that one can infer key-pressing actions by considering only lower-arm muscles. To design such a system, we require choosing (i) the type of sensor, (ii) the input it provides to the inference algorithm, and (iii) the inference algorithm itself. For our system, we chose to use an US sensor capable of recording images at a high frequency, to capture muscle structure. Choosing an US relies on previous work [12] in

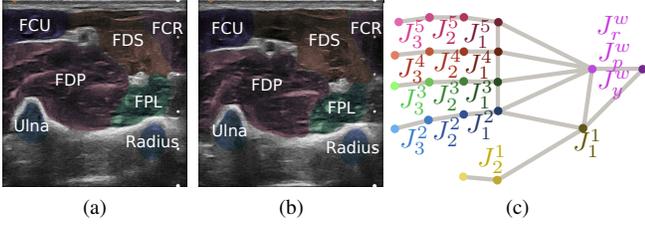


Fig. 2: (a,b) Visual segmentation of two US images (top of each image is the skin shell). The Ulna and Radius (the two main bones), are marked in blue, and the flexion muscles are marked with different colors. All fingers are at rest in the left image while in the right image the Ring finger is fully flexed. Notice that the region capturing the FDS in the right image is larger than on the left image as the FDS becomes more dominant in this area during movement. (c) Visual illustration of a hand's configuration. The Thumb (F^1) is associated with two joints [J_1^1, J_2^1], the other four fingers F^i ($i \in \{2, \dots, 5\}$) are each associated with three joints [J_1^i, J_2^i, J_3^i] and three additional DOFs are associated with the wrist [J_r^w, J_p^w, J_y^w]. The set of joints $\{J_1^i | 1 \leq i \leq 5\}$, $\{J_2^i | 2 \leq i \leq 5\}$ and $\{J_3^i | 2 \leq i \leq 5\}$ correspond to the MP, PIP and DIP joints, respectively (see Sec. III). Figures are best viewed in color.

which US completely outperformed EMG in classifying gestures. Moreover, when compared to the single-dimensional signal provided by an EMG sensor, spatial features that are found in an US image have more information about fingers' whereabouts, and are less noisy. This will be key in the way we address the latter two design choices.

Recall that an US provides us with an image. However, in contrast to existing approaches [5], [16], we chose to use a *sequence* of the latest-acquired US images as the input to our inference algorithm. As we will demonstrate empirically (Sec. VI), considering a sequence of US images allows to exploit the dynamics which in turn, leads to significantly better results in downstream tasks.

Given a stream of US images we chose to solve the inference problem via a NN architecture that is tailored for processing images, exploits temporal coherence, and allows a domain expert to inject domain knowledge via an intermediate representation. In our setting, this intermediate representation is an encoding of the hand modeled as a robotic manipulator. We detail the exact architecture and specific design considerations in Sec. VI.

V. INFERENCE LEARNING FRAMEWORK

In this section, we describe our inference learning framework. We start (Sec. V-A) by formally defining and modeling the inference problem and then continue (Sec. V-B) to describe an approach where we only make use of US images without the aforementioned intermediate representation. As we will see (Sec. VI), utilizing an intermediate representation during training can lead to better inference capabilities. Thus, in Sec. V-C we extend our first method to account for said representation. Notice that in this section we concentrate on the high-level learning framework and the exact training and inference procedures are detailed in Sec. VI-B.

A. Notation and Model

In this section, we introduce the necessary notation and formally define our inference problem. We assume to be given a sequence of k US images corresponding to k timestamps and wish to use them to predict which finger was pressed for each timestamp k .

Let F^1, \dots, F^5 denote the five fingers with F^1 and F^5 corresponding to the Thumb and Little finger, respectively. For each finger F^i and each timestamp t_j , we associate a binary variable $p_j^i \in \{0, 1\}$ such that $p_j^i = 1$ if finger F^i is pressed at timestamp t_j and $p_j^i = 0$ otherwise. For every timestamp t_j , we denote by $p_j \in \{0, 1\}^5$ the binary vector indicating which fingers are pressed at t_j which we term a *pressing vector*. Namely, $p_j := \langle p_j^1, \dots, p_j^5 \rangle$. We assume that for every timestamp t_j there is a corresponding US image s_j and *configuration* x_j . Here, x_j is an encoding of the fingers and wrists as a set of joint angles (see Fig. 2c).

Our training data consists of tuples $\langle \bar{s}_\ell, \bar{x}_\ell, \bar{p}_\ell \rangle$. Here, $\bar{s}_\ell := \langle s_\ell, \dots, s_{k+\ell} \rangle$ is a sequence of k consecutive US images, $\bar{x}_\ell := \langle x_\ell, \dots, x_{k+\ell} \rangle$ is a sequence of k consecutive configurations and $\bar{p}_\ell := \langle p_\ell, \dots, p_{k+\ell} \rangle$ is a sequence of k consecutive pressing vectors all of which start at timestamp t_ℓ . We denote by \mathcal{S} and \mathcal{X} the sets of all sequences of k US images and configurations, respectively. For each finger F^i and each timestamp t_j we wish to predict p_j^i (i.e., if F^i was pressed at t_j). To this end, our system will output the variable $\hat{p}_j^i \in [0, 1]$ estimating the probability that $p_j^i = 1$. Denote $\hat{p}_j := \langle \hat{p}_j^1, \dots, \hat{p}_j^5 \rangle$ and $\hat{p}_\ell := \langle \hat{p}_\ell, \dots, \hat{p}_{k+\ell} \rangle$. As we will see, our approach will make use of an intermediate representation corresponding to the configuration at each timestamp. Thus, we will also denote \hat{x}_j as the system's estimation of x_j and $\hat{x}_\ell := \langle \hat{x}_\ell, \dots, \hat{x}_{k+\ell} \rangle$. In Sec. VI we describe the experimental setting which allows us to obtain labeled training data for our representative tasks. In the meantime, we describe our approach to solving our inference problem.

B. Learning Pressing Events Directly

To estimate pressing vectors, we use a hybrid architecture starting with a convolutional neural network (CNN) module [25] to extract spatial features from the US images, followed by a recurrent neural network (RNN) [26], [27] to aggregate these features and infer the pressing vectors. Such a combination was shown to be effective in various computer-vision tasks that require temporal coherence [28], [29], [30]. By combining these modules, we create a sequence-to-sequence model, which receives a sequence of US images \bar{s}_ℓ as an input and outputs \hat{p}_ℓ which estimate the pressing vector \bar{p}_ℓ . We call this the multi-frame (MF) model and denote it as G_θ where θ is the set of model parameters that we optimize in the training phase. This is done by minimizing the negative log-likelihood, i.e.,

$$\mathcal{L}_G(\theta) := -\mathbb{E}_{\bar{s}_\ell \sim \mathcal{S}} [\log P(\hat{p}_\ell | \bar{s}_\ell, \theta)]. \quad (1)$$

C. Learning Pressing Events via Hand Configurations

To estimate pressing vectors while making use of an intermediate representation of the finger's configuration, we

use an auto-encoder [26], [31] framework composed of two main blocks—an encoder mapping the sequence of US images \bar{s}_ℓ to a sequence of configurations \hat{x}_ℓ and a decoder using \hat{x}_ℓ to estimate the pressing vector by outputting \hat{p}_ℓ . The encoder’s architecture follows the architecture proposed for the MF model (i.e., a CNN followed by an RNN) while the decoder uses an additional RNN module together with a residual connection [32] that feeds the model with the hidden features from the encoder (see Fig. 1 for visualization).

Auto-encoders allow a domain expert to inject domain knowledge via an intermediate representation (hand configurations in our setting) [33], [34]. The residual connection was added to better allow the decoder to infer pressing vectors by obtaining knowledge from both US images and configurations. We name this model the Configuration-Based Multi-Frame (CBMF) model, and denote it as $H_{\theta,\phi} = h_\phi^d \circ h_\theta^e$ where θ and ϕ are the encoder’s and decoder’s parameters, respectively which we optimize in the training phase.²

While the exact details of how we train this model are explained in the experimental section (Sec. VI-B), we note that the training process is done in two phases. In the first phase, we train only the encoder h_θ^e while in the second phase, we simultaneously train the (pre-trained) encoder h_θ^e together with the (untrained) decoder h_ϕ^d . Splitting the learning into two phases (i.e., a single and a multi-modal supervised learning task) was shown to produce better results in similar tasks [35], [36], [37]. Roughly speaking, by pre-training h_θ^e we get a “coarse” mapping between US images \bar{s}_ℓ and predicted hand configurations \hat{x}_ℓ . This allows us to “bootstrap” the network’s parameters when solving the entire inference problem in our multi-modal task. More formally, given \bar{s}_ℓ , \hat{x}_ℓ , and \hat{p}_ℓ , the auto-encoder is trained to predict \hat{x}_ℓ and \hat{p}_ℓ given \bar{s}_ℓ by minimizing the negative log-likelihood of the combination of both terms, i.e.,

$$\mathcal{L}_H(\theta, \phi) = -\mathbb{E}_{\bar{s}_\ell \sim \mathcal{S}, \hat{x}_\ell \sim \mathcal{X}} [\log P(\hat{x}_\ell | \bar{s}_\ell, \theta) + \log P(\hat{p}_\ell | \hat{x}_\ell, \bar{s}_\ell, \phi)]. \quad (2)$$

VI. EXPERIMENTS

In this section, we evaluate our approach for predicting fine finger motions. We start (Sec. VI-A) by describing our experimental setup and data-collection methodology. We then describe our model and training details (Sec. VI-B) and finish with an analysis of the data collected (Sec. VI-C).

A. Experimental Setup and Data Collection

Our setup includes a Clarius L15HD (wireless B-mode US device), a Roland FP-30 (digital piano), a standard keyboard, and a Vicon motion-capture system. The US’s configuration was set for the musculoskeletal task (MSK), with 10 MHz frequency generating a 5 cm \times 5 cm processed image at a rate of 19-21 images per second, in a resolution of 480 \times 480. It was attached using a wearable strap (see Fig. 1) to the so-called “transverse” location (see, e.g., [9]). The piano and keyboard were used to record the set of pressed notes and

²Here the superscript ‘e’ and ‘d’ are used to refer to the encoder and decoder modules, respectively.

keys during each US frame and the Vicon system was used to track the palm’s joint locations. Here, each joint location is captured and used to extract the joint angles that form the hand’s configuration, as described in the appendix (Sec. VII).

We collected data from twelve subjects; averaging at 23 years old, seven of whom are men and five of whom are women. All subjects were confirmed to be right-handed, experienced in playing piano, and without any neurological disorders or deformity in the relevant hand area. At the beginning of both tasks, subjects were instructed to put their hands in a stationary state in which each of their fingers must be on the same note/key for the entire session. For piano playing, subjects were asked to play melodies that only require five distinguished notes and can be played using one hand only. For keyboard typing, subjects were asked to randomly type on each of the five keys their fingers are stationed over, while not using more than one key at the same time. For each task, subjects were asked to perform 2-3 sessions, with approximately ninety seconds for each session, and rest between sessions. The process was approved by the institution’s Ethics Committee. The final dataset includes 44,596 samples for piano playing and 42,143 samples for keyboard typing, each sample representing the tuple $\langle \bar{s}_\ell, \hat{x}_\ell, \hat{p}_\ell \rangle$ defined in Sec. V-A.

B. Model Training

Before training, both images s_i and configurations x_i were normalized to the range of [0, 1]. On feedforward, our MF model G_θ receives a sequence of $k = 8$ US images \bar{s}_ℓ , such that each image s_i is downsampled to a resolution of 224 \times 224 and fed to the CNN block to extract feature maps, which are then flattened, concatenated and fed into the RNN module to output k pressing vectors. Our CBMF model $H_{\theta,\phi}$ works such that the encoder h_θ^e is being fed with the images similar to the MF model, and outputs k configurations, that are then fed to the decoder h_ϕ^d (an additional RNN block), along with the features from the residual connection. The decoder h_ϕ^d is designed to output k predicted configurations. The entire architecture is detailed in the appendix (Sec. VII). On back propagation, for the predicted configurations, the applied loss function for $N = 32$ samples in a batch is the Mean Squared Error (MSE) on all configurations, i.e., $\mathcal{L}_{\text{MSE}} = \frac{1}{N \cdot k} \sum_{i \in N \cdot k} \|\hat{x}_i - x_i\|_2^2$. A configuration x_i is a vector of 17 angles, as depicted in Fig. 2c. For the predicted pressing vectors, the applied loss objective for $N = 32$ samples in a batch, such that each sample is k vectors of $M = 5$ probabilities, is the Binary Cross Entropy (BCE) for all probabilities, i.e., $\mathcal{L}_{\text{BCE}} = \frac{1}{N \cdot k \cdot M} \sum_{i \in M} \sum_{j \in N \cdot k} p_j^i \log \hat{p}_j^i + (1 - p_j^i) \log(1 - \hat{p}_j^i)$. To train the MF model, we apply the loss function \mathcal{L}_{BCE} , and to train the CBMF model, we apply the combined loss function $\lambda \mathcal{L}_{\text{MSE}} + \mathcal{L}_{\text{BCE}}$, where λ is chosen using a trial-and-error procedure (see the the appendix (Sec. VII)). All training sessions were executed for 20 epochs, to avoid overfitting. Adam [38] was used as the optimization method, with a learning rate of 0.001.

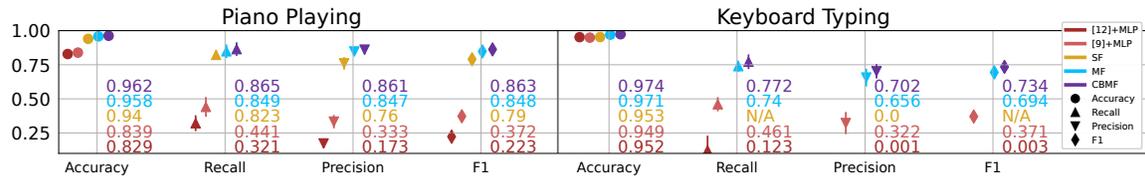


Fig. 3: Evaluation of the MF and CBMF methods (Sec. V-B and V-C) as well as the three baselines (Sec. VI-C) for piano playing and keyboard typing tasks. Accuracy, recall, precision, and F1 were averaged over all test samples in 5 folds. The shape in each interval presents the mean values, while the vertical segment represents the standard deviation range.

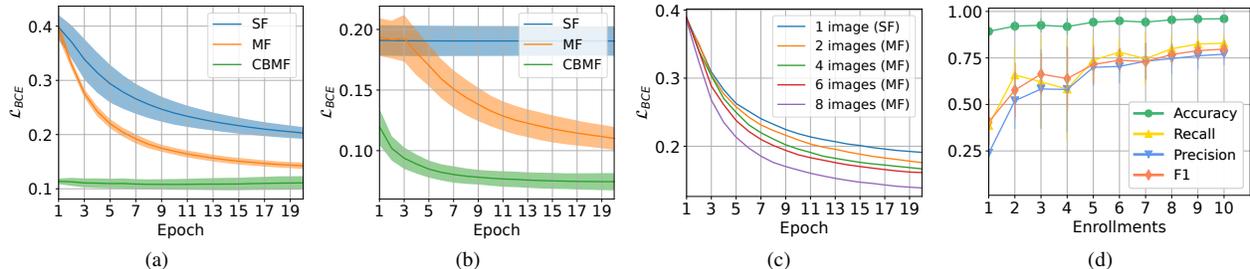


Fig. 4: Comparison of the optimization process using \mathcal{L}_{BCE} as a function of the epoch number for the piano-playing (a) and keyboard-typing (b) tasks. (c) Training loss as a function of epoch number for different sizes of US sequences as inputs. (d) Evaluation of the CBMF model on piano playing as a function of the number of enrollments. At each step, we randomly add a new enrollment to the training set and evaluate an unseen enrollment. Results are averaged over 10 experiments.

C. Results

In all of the experiments, we evaluated all methods using 5-fold validation: for each task, all sessions of all subjects were divided into five folds, and trained such that the i 'th fold was set aside for testing and not used in training. This ensures that the test set contains unseen intervals of motions.

How good are the proposed methods? We evaluated the inference capabilities of our two models (MF and CBMF) on the piano-playing and keyboard-typing tasks by plotting the accuracy, recall, precision, and $F1$ (Fig. 3) as well as the test loss during training (Fig. 4a and 4b). For baselines, we added the single-frame (SF) model, where we simply take our MF model and replace the RNN with a multilayer perceptron (MLP) to predict a single pressing vector from a single image. In addition, we implemented the state-of-the-art algorithms for gesture classification [9], [12]. To allow the detection of pressing vectors, we concatenated an MLP to their output. Implementation details are in the appendix (Sec. VII). We also tested common computer-vision models, such as the VGG16 [39], [16] model (pretrained on ImageNet [40], or not), and C3D [41], which were shown to perform favorably in video classification. However, the models did not converge in any of the tasks and therefore were excluded from evaluation.

The overall performance for piano playing is significantly higher than for keyboard typing. We associate this difference with differences in pressing duration and body postures between the tasks: The maximum note offset on our piano and our keyboard is 10.45 mm and 2.95 mm, respectively. This implies longer muscle flexing when playing the piano, which leads to more dominant motions in US images. Additionally, playing the piano involves different body postures which often lead to more dominant motions in US images.

Moving on to comparing the different models, the SF

model is unable to recreate the finer task of keyboard typing, but is able to learn to play the piano. This backs up our conjecture that a sequence of multiple frames is required to infer fine finger motions. Notice that the MF model outperforms the three baselines while the CBMF model outperforms all. This is especially noticeable in the harder keyboard-typing task (see recall and precision), which backs up our conjecture that our intermediate representation allows us to better infer fine finger motions, by reducing false positive and false negative predictions. In addition, we examined piano playing with sequences of images of different sizes ($k \in \{1, 2, 4, 6, 8\}$), showing that convergence on the test set improves gradually the more we increase k . Larger sequences may have been beneficial, however, we did not use more than 8 consecutive images due to GPU memory limitations. For a representative visualization of the difference between the MF and CBMF models, see Fig. 5.

How feasible is the method for a group of subjects? We plot the different metrics for each subject individually to see how the CBMF model varies across different subjects (Fig. 6). Overall, the results for piano playing are better than for keyboard typing, and those who performed better on piano playing, performed better on keyboard typing, but most importantly, the method works for all subjects, showing us the potential for large-scale adoption of these algorithms.

How well does the system work for all fingers? We computed the recall and precision values for each of the fingers individually to understand how each of the fingers affects the quality of the solution (Fig. 6). On both tasks, the model performs better on fingers F^2 , F^3 and F^4 , rather than on fingers F^1 and F^5 . A possible explanation is that operating the three middle fingers causes larger movement in the muscles that are found in the US images when compared to the Thumb and Little finger. As illustrated in Fig. 2, the muscle responsible for Thumb movement is significantly less

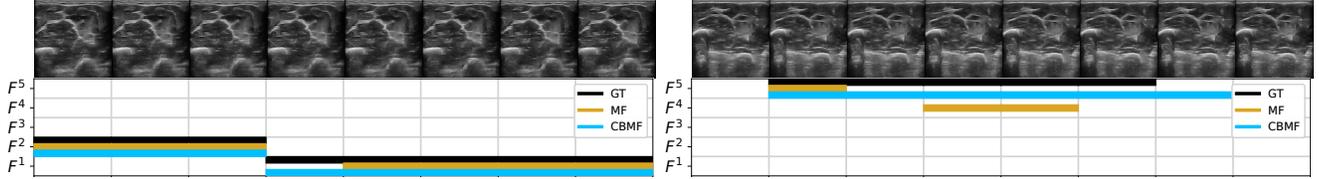


Fig. 5: Samples from executing the MF and CBMF models. In each rollout, the upper row represents finger F^5 (Little finger) and the lower row represents finger F^1 (Thumb). Ground-truth (GT) labels are marked in black. The differences between both models are expressed in consistency and the number of false-positive predictions. Notice from the US images how hard is it for a human eye to understand subtle morphological changes during such gentle movements.

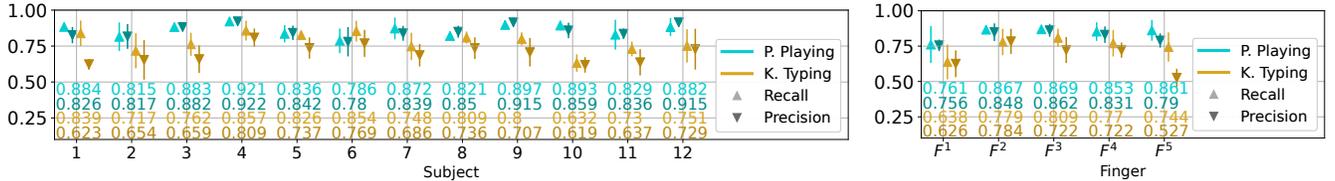


Fig. 6: (Left) Recall and precision for the CBMF model for all twelve subjects. Vertical segments denote one standard deviation. (Right) Recall and precision separated for all five fingers, showing the expected performance of each finger.

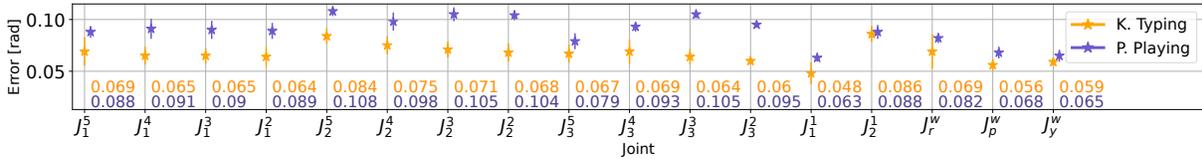


Fig. 7: Error-values in radians for predicting joint angles using our CBMF model.

dominant than the rest of the muscles, which might explain why Thumb movement is the least-understandable gesture.

How well can the model generalize on real-time data?

Our method works without any preprocessing, and can work in an online setting. This is in contrast to existing methods for which subjects were instructed to perform gestures according to predefined instructions [5], or separate between gesture categories before feature extraction [9].

A potential challenge in real-time inference is that the US sensor is not placed in the exact location for which training data was obtained (also known as sensor reattachment). To this end, we define “enrollment” as a period of 2-3 recording sessions, in which the US sensor was not removed from the hand. We used a single subject to demonstrate the number of enrollments needed to infer from an unseen enrollment. Results (Fig. 4d) show that (i) as the number of enrolments increases in the training session, our inference capabilities increase and (ii) the number of enrolments required in the training session is relatively moderate. We connected our system to a robotic hand (see Fig. 8 as well as the supplementary video) which continuously replayed in real time notes played on the piano.



Fig. 8: Our system connected to a robotic hand.

How accurate are the estimated configurations? We evaluated the error values (in radians) of our CBMF model for each of the individual 17 joints in a configuration (see Fig. 2c). Results, depicted in Fig. 7, demonstrate that the more the task requires longer movements, the higher the average error, especially for the more dominant joints, such as fingers F^2 and F^3 . Interestingly, the DOFs that are more dominant near the wrist ($J_1^1, J_r^w, J_p^w, J_y^w$), achieve a smaller error on average, which might relate to the large region captured by the relative muscles in the US images.

VII. DISCUSSION AND FUTURE WORK

In this paper, we explored the task of predicting and differentiating between fine finger motions by only having access to a stream of lower-arm US images. We show that incorporating temporal coherence and a kinematic hand representation in data-driven methods is essential for solving such tasks. We believe that the inference capabilities we present will pave the way towards enabling fully-operational prostheses for those in need. Importantly, generalization for unseen subjects was not dealt with in this paper, as it requires holding prior knowledge regarding the muscles’ structure.

Future directions that we are exploring include (i) Restoring hand motions for people with muscular deformations such as amputees. This will potentially require data-driven solutions to learn continuous and unlabeled finger motions for non-able-bodied subjects. (ii) Recent advancements in feedback sensing allow amputees to sense objects they are engaging with [42]. Sensor feedback can also be used to endow our inference algorithm with additional input such as surface type which can be used to improve decision making.

ACKNOWLEDGMENT

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 863839), the Israeli Ministry of Science & Technology grants No. 3-16079 and 3-17385, and the United States-Israel Binational Science Foundation (BSF) grants no. 2019703 and 2021643. We also thank Haifa3D, a non-profit organization providing self-made 3d-printed solutions, for their consulting and support through the research.

APPENDIX

Joints Computation

The set of joint angles forming a configuration is acquired by placing markers on the relevant joints and using a motion-capture system to obtain the $[x, y, z]$ location of each marker (see Fig. 9). Then, each joint angle is computed using the markers’ locations as we detail next. Hereafter, we compute vectors as the point-wise subtraction of the marker locations. For example, given two 3D points (that can be considered as two 3D vectors w.r.t. some reference frame) f_{42} and f_{43} , we denote the vector from f_{42} to f_{43} as follows:

$$\vec{v}_{f_{42} \rightarrow f_{43}} = \vec{f}_{43} - \vec{f}_{42}. \quad (3)$$

Next, we describe how we compute the angle between two vectors. Consider, for example, the vector from f_{42} to f_{43} and the vector from f_{42} to f_{41} , we compute the joint angle J_2^5 as follows:

$$J_2^5 = \arctan2(\|\vec{v}_{f_{42} \rightarrow f_{41}} \times \vec{v}_{f_{42} \rightarrow f_{43}}\|_2, \vec{v}_{f_{42} \rightarrow f_{41}} \cdot \vec{v}_{f_{42} \rightarrow f_{43}}). \quad (4)$$

This allows us to compute joint angles. We start with the four fingers F^2, \dots, F^5 and will use finger F^5 (Little finger) for example. The three joints that are required are J_1^5 , J_2^5 and J_3^5 . Hence, we will compute J_1^5 and J_3^5 as follows (fingers F^1 , F^2 , and F^3 are calculated in a similar fashion):

$$\begin{aligned} J_1^5 &= \arctan2(\|\vec{v}_{f_{41} \rightarrow t_5} \times \vec{v}_{f_{41} \rightarrow f_{42}}\|_2, \\ &\quad \vec{v}_{f_{41} \rightarrow t_5} \cdot \vec{v}_{f_{41} \rightarrow f_{42}}), \\ J_3^5 &= \arctan2(\|\vec{v}_{f_{42} \rightarrow f_{41}} \times \vec{v}_{f_{42} \rightarrow f_{43}}\|_2, \\ &\quad \vec{v}_{f_{42} \rightarrow f_{41}} \cdot \vec{v}_{f_{42} \rightarrow f_{43}}). \end{aligned} \quad (5)$$

As for the Thumb, we compute two joint angles; J_1^1 and J_2^1 . Similar to Eq. 5, we compute J_2^1 as follows:

$$J_2^1 = \arctan2(\|\vec{v}_{t_6 \rightarrow t_7} \times \vec{v}_{t_6 \rightarrow t_5}\|_2, \vec{v}_{t_6 \rightarrow t_7} \cdot \vec{v}_{t_6 \rightarrow t_5}), \quad (6)$$

where joint J_1^1 is acquired by computing the angle between $\vec{v}_{t_5 \rightarrow t_6}$ and the normal vector to the plane approximating the palm’s plane, by taking the three points f_{11} , f_{41} and t_4 . Denoting this normal as \hat{n}_p , joint J_1^1 is computed as follows:

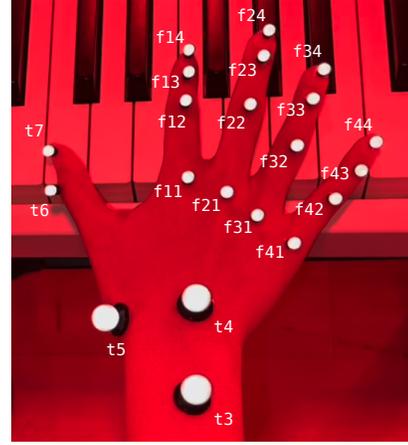


Fig. 9: An image showing the position of each marker on the subject’s hand.

$$J_1^1 = \arctan2(\|\vec{v}_{t_5 \rightarrow t_6} \times \hat{n}_p\|_2, \vec{v}_{t_5 \rightarrow t_6} \cdot \hat{n}_p). \quad (7)$$

This is done differently in order to compute the angle that corresponds to the pressing motion of the Thumb. As for the three DOFs of the wrist, we have J_r^w , J_p^w , and J_y^w for roll, pitch and yaw respectively. As for pitch, we compute the angle between \hat{n}_p and $\vec{v}_{t_4 \rightarrow t_3}$. For yaw, we compute the angle between $\vec{v}_{f_{41} \rightarrow f_{11}}$ and $\vec{v}_{t_4 \rightarrow t_3}$, and for roll, we compute the angle between \hat{n}_p and a vector obtained from two additional points located on the elbow.

Network Architecture

Tables I and Table II shows the network architectures for the MF and CBMF models, respectively. The tables are divided into blocks, such that the first block represents the CNN module and the others represent the RNN modules. The input shape for both models is $[224 \times 224 \times 8]$, and corresponding to 8 consecutive grayscale US images. The shape of the intermediate representation is $[8 \times 17]$, corresponding to a trajectory of 8 robotic configurations, each of size 17. The output shape is $[8 \times 5]$, denoting 8 pressing vectors, each containing five probabilities for the five fingers.

The CNN module is based on the encoding part of UNet [43], due to the high success of this architecture when used in medical imaging tasks. This block is used to extract 512 3×3 feature maps from each of the eight images, such that the same parameters are shared across all images. These are flattened, concatenated, and fed into the recurrent block. In the CBMF model, The trajectory, along with 8×128 features are taken from the output of the second GRU layer, concatenated, and fed as an input to the decoder.

Baselines Implementation

As presented in Sec. VI-C, we implemented three baselines to assess our models in solving the two downstream tasks. The evaluation was using the same 5-fold validation over the same datasets. We started by evaluating the importance of temporal coherence by defining the SF model.

Layer	Act.	Out. shape	Parameters
Conv2D	ReLU	224 × 224 × 32	320
Conv2D	ReLU	224 × 224 × 32	9,248
MaxPooling2D (2 × 2)	-	112 × 112 × 32	-
Conv2D	ReLU	112 × 112 × 64	18.50k
Conv2D	ReLU	112 × 112 × 64	36.93k
MaxPooling2D (2 × 2)	-	56 × 56 × 64	-
Conv2D	ReLU	56 × 56 × 128	73.86k
Conv2D	ReLU	56 × 56 × 128	147.58k
MaxPooling2D (2 × 2)	-	28 × 28 × 128	-
Conv2D	ReLU	28 × 28 × 256	295.17k
Conv2D	ReLU	28 × 28 × 256	590.08k
Dropout (p=0.3)	-	28 × 28 × 256	-
MaxPooling2D (2 × 2)	-	14 × 14 × 256	-
Conv2D	ReLU	14 × 14 × 512	1.18M
Conv2D	ReLU	14 × 14 × 512	2.36M
Dropout (p=0.3)	-	14 × 14 × 512	-
MaxPooling2D (4 × 4)	-	3 × 3 × 512	-
Flatten	-	4608	-
<hr/>			
Concatenate (8 images)	-	8 × 4608	-
GRU	ReLU	8 × 1024	17.31M
GRU	ReLU	8 × 128	443.14K
GRU	Linear	8 × 5	2,025
<hr/>			
Total			22.93M

TABLE I: Neural-network architecture of our MF model G_θ .

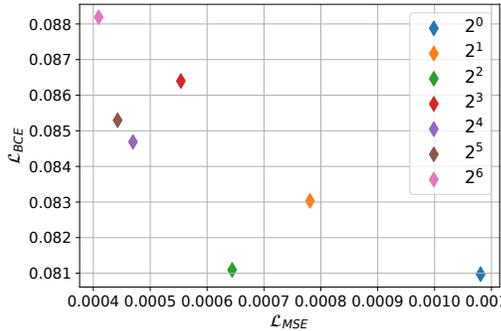


Fig. 10: \mathcal{L}_{BCE} and \mathcal{L}_{MSE} for different values of λ . On average, results are showing that using lower values of λ yields lower values of \mathcal{L}_{BCE} and higher values of λ yields lower values of \mathcal{L}_{MSE} and higher values of \mathcal{L}_{BCE} .

Here, we take our MF model and replace the RNN with a MLP, such that the CNN module is the same as in the MF and CBMF models, and the addition is two linear layers. On inference, the model receives a single 224×224 US image and predict a single pressing vector.

In addition we compared our work to two additional papers that succeeded in gesture classification. Importantly, these papers were originally designed to classify full recorded sessions in an offline setting, and since these algorithms were evaluated on different gestures with different sensors, we modified the input to receive the last acquired 8 images as in our models. Second, these algorithms were designed to discriminate between gestures and here we allow multiple fingers pressing at the same time. Therefore, we set the amount of available classes to $2^5 = 32$ (i.e., a class corresponds to a subset of fingers pressed).

Huang et al. [12] proposed to sample five columns of pixels from the image, imitating the feature structure of A-

Layer	Act.	Out. shape	Parameters
Conv2D	ReLU	224 × 224 × 32	320
Conv2D	ReLU	224 × 224 × 32	9,248
MaxPooling2D (2 × 2)	-	112 × 112 × 32	-
Conv2D	ReLU	112 × 112 × 64	18.50k
Conv2D	ReLU	112 × 112 × 64	36.93k
MaxPooling2D (2 × 2)	-	56 × 56 × 64	-
Conv2D	ReLU	56 × 56 × 128	73.86k
Conv2D	ReLU	56 × 56 × 128	147.58k
MaxPooling2D (2 × 2)	-	28 × 28 × 128	-
Conv2D	ReLU	28 × 28 × 256	295.17k
Conv2D	ReLU	28 × 28 × 256	590.08k
Dropout (p=0.3)	-	28 × 28 × 256	-
MaxPooling2D (2 × 2)	-	14 × 14 × 256	-
Conv2D	ReLU	14 × 14 × 512	1.18M
Conv2D	ReLU	14 × 14 × 512	2.36M
Dropout (p=0.3)	-	14 × 14 × 512	-
MaxPooling2D (4 × 4)	-	3 × 3 × 512	-
Flatten	-	4608	-
<hr/>			
Concatenate (8 images)	-	8 × 4608	-
GRU	ReLU	8 × 1024	17.31M
GRU	ReLU	8 × 128	443.14K
GRU	Linear	8 × 17	7,497
<hr/>			
Merge	-	8 × 145	-
GRU	ReLU	8 × 256	309.50K
GRU	ReLU	8 × 128	148.22K
GRU	Linear	8 × 5	2,025
<hr/>			
Total			22.93M

TABLE II: Neural-network architecture of our CBMF model $H_{\theta,\phi}$.

mode transducers. These columns are then separated into segments, from which they extracted two coefficients using linear fitting. Then, these features are concatenated across the temporal axis and fed to the classifier. Here, the authors used sequences of 36 images, but since we did not observe any improvement, the results are presented for sequences of eight images. Next, McIntosh et al. [9] proposed to extract the optical flow from adjacent frames, perform average pooling on patches of size 20×20 , and accumulate them across the temporal axis. In the same way, we calculated these features for a sequence of eight US images and fed them to the classifier. For both baselines, we evaluated all classifiers that were proposed by the authors, and in both cases, MLP outperformed all other classifiers.

Losses Weighting

We empirically chose the value of λ to minimize the objective function $\lambda\mathcal{L}_{MSE} + \mathcal{L}_{BCE}$. Here, we executed training sessions to obtain results for $\lambda \in \{2^i | i \in [0, \dots, 6]\}$ and chose the result yielding the lowest \mathcal{L}_{BCE} , without harming the optimization of \mathcal{L}_{MSE} . Results, presented in Fig. 10, show that by taking $\lambda = 4$, we benefit from using both modalities. The same experiment with $\lambda = 0$ (i.e., only using \mathcal{L}_{BCE}) resulted in $\mathcal{L}_{BCE} = 0.093$ and $\mathcal{L}_{MSE} = 2.582$, and therefore was ruled out from the evaluation.

Real-Time Implementation

We built a setup for online real-time inference to demonstrate piano playing on demand. The code is written in C++ and includes multi-threading to deal with interfaces to the US

sensor, the robotic hand, and model inference simultaneously. For optimized inference, our CBMF model was divided into two software modules - the encoder's convolution block and the remaining recurrent blocks including the residual connection. On inference, the incoming images from the US sensor are fed to the first module to extract features, while an additional thread, simultaneously feeds the sequence of the last eight extracted features to the second module, and stores the predictions as instructions that are sent to the robotic hand. The model was optimized using the ONNX Runtime library, the entire model weighs 745 MB and the inference is executed at 60 Hz on a powerful GPU (Nvidia RTX 2080Ti) or 29 Hz on a low-budget GPU (Nvidia Quadro P620).

REFERENCES

- [1] C. Connolly, "Prosthetic hands from touch bionics," *Ind. Rob.*, 2008.
- [2] P. Herberts, "Myoelectric signals in control of prostheses: Studies on arm amputees and normal individuals," *Acta. Orthop.*, vol. 40, no. sup124, pp. 1–83, 1969.
- [3] J. A. George, T. S. Davis, M. R. Brinton, and G. A. Clark, "Intuitive neuromyoelectric control of a dexterous bionic arm using a modified kalman filter," *arXiv*, vol. abs/1908.10522, 2019.
- [4] S. Said, M. Sheikh, F. Al-Rashidi, Y. Lakys, T. Beyrouthy, A. Nait-ali *et al.*, "A customizable wearable robust 3d printed bionic arm: Muscle controlled," in *BioSMART*, 2019, pp. 1–6.
- [5] A. S. Dhawan, B. Mukherjee, S. Patwardhan, N. Akhlaghi, G. Levay, R. Holley, W. M. Joiner, M. Harris-Love, and S. Sikdar, "Proprioceptive somyographic control: A novel method of intuitive proportional control of multiple degrees of freedom for upper-extremity amputees," *arXiv*, vol. abs/1808.06543, 2018.
- [6] J. Yan, X. Yang, X. Sun, Z. Chen, and H. Liu, "A lightweight ultrasound probe for wearable human-machine interfaces," *IEEE Sens. J.*, vol. 19, no. 14, pp. 5895–5903, 2019.
- [7] X. Yang, D. Zhou, Y. Zhou, Y. Huang, and H. Liu, "Towards zero re-training for long-term hand gesture recognition via ultrasound sensing," *IEEE J. Biomed. Health Informatics*, vol. 23, no. 4, pp. 1639–1646, 2019.
- [8] X. Yang, J. Yan, Y. Fang, D. Zhou, and H. Liu, "Simultaneous prediction of wrist/hand motion via wearable ultrasound sensing," *IEEE Trans. Neural Syst. Rehabilitation Eng.*, vol. 28, no. 4, pp. 970–977, 2020.
- [9] J. McIntosh, A. Marzo, M. Fraser, and C. Phillips, "Echoflex: Hand gesture recognition using ultrasound imaging," in *CHI*, 2017, pp. 1923–1934.
- [10] L. Resnik, M. R. Meucci, S. Lieberman-Klinger, C. Fantini, D. L. Keltly, R. Disla, and N. Sasson, "Advanced upper limb prosthetic devices: implications for upper limb prosthetic rehabilitation," *Arch. Phys. Med. Rehabil.*, vol. 93, no. 4, pp. 710–717, 2012.
- [11] U. Côté-Allard, C. L. Fall, A. Drouin, A. Campeau-Lecours, C. Gosselin, K. Glette, F. Lavolette, and B. Gosselin, "Deep learning for electromyographic hand gesture signal classification by leveraging transfer learning," *arXiv*, vol. abs/1801.07756, 2018.
- [12] Y. Huang and H. Liu, "Performances of surface EMG and ultrasound signals in recognizing finger motion," in *HSI*, 2016, pp. 117–122.
- [13] Y. Herbst, L. Zelnik-Manor, and A. Wolf, "Analysis of subject specific grasping patterns," *PLoS One*, vol. 15, no. 7, p. e0234969, 2020.
- [14] K. Bimbraw, E. Fox, G. Weinberg, and F. L. Hammond, "Towards somyography-based real-time control of powered prosthesis grasp synergies," in *EMBC*, 2020, pp. 4753–4757.
- [15] X. Yang, Z. Chen, N. Hettiarachchi, J. Yan, and H. Liu, "A wearable ultrasound system for sensing muscular morphological deformations," *IEEE Trans. Syst. Man Cybern. Syst.*, vol. 51, no. 6, pp. 3370–3379, 2021.
- [16] K. Bimbraw, C. J. Nycz, M. Schueler, Z. Zhang, and H. K. Zhang, "Prediction of metacarpophalangeal joint angles and classification of hand configurations based on ultrasound imaging of the forearm," *arXiv*, vol. abs/2109.11093, 2021.
- [17] Y. Wang, X. Ge, H. Ma, S. Qi, G. Zhang, and Y. Yao, "Deep learning in medical ultrasound image analysis: A review," *IEEE Access*, vol. 9, pp. 54 310–54 324, 2021.
- [18] S. Vedula, O. Senouf, A. M. Bronstein, O. V. Michailovich, and M. Zibulevsky, "Towards ct-quality ultrasound imaging using deep learning," *arXiv*, vol. abs/1710.06304, 2017.
- [19] R. J. G. van Sloun, R. Cohen, and Y. C. Eldar, "Deep learning in ultrasound imaging," *arXiv*, vol. abs/1907.02994, 2019.
- [20] S. Vedula, O. Senouf, G. Zurakhov, A. M. Bronstein, O. V. Michailovich, and M. Zibulevsky, "Learning beamforming in ultrasound imaging," in *MIDL*, vol. 102, 2019, pp. 493–511.
- [21] H. Bliss, S. Bird, A. P. Cooper, S. Burton, and B. Gick, "Seeing speech: Ultrasound-based multimedia resources for pronunciation learning in indigenous languages," 2018.
- [22] N. Kimura, M. Kono, and J. Rekimoto, "Sottovoce: An ultrasound imaging-based silent speech interaction using deep neural networks," in *CHI*, 2019, p. 146.
- [23] F. Vogt, G. McCaig, M. A. Ali, and S. S. Fels, "Tongue 'n' groove: An ultrasound based music controller," in *NIME*, 2002, pp. 60–64.
- [24] S. M. Lulich, S. Charles, and B. Lulich, "The relation between tongue shape and pitch in clarinet playing using ultrasound measurements," *J. Acoust. Soc. Am.*, vol. 141, no. 3, pp. 1759–1768, 2017.
- [25] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroury, B. Shuai, T. Liu, X. Wang, and G. Wang, "Recent advances in convolutional neural networks," *arXiv*, vol. abs/1512.07108, 2015.
- [26] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," Tech. Rep., 1985.
- [27] A. Sherstinsky, "Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network," *arXiv*, vol. abs/1808.03314, 2018.
- [28] J. Guo, Y. Liu, Q. Yang, Y. Wang, and S. Fang, "Gps-based citywide traffic congestion forecasting using cnn-rnn and c3d hybrid model," *Transportmetrica A*, vol. 17, no. 2, pp. 190–211, 2021.
- [29] W. Lu, J. Li, Y. Li, A. Sun, and J. Wang, "A cnn-lstm-based model to forecast stock prices," *Complex.*, vol. 2020, pp. 6622927:1–6622927:10, 2020.
- [30] Q. Zou, H. Jiang, Q. Dai, Y. Yue, L. Chen, and Q. Wang, "Robust lane detection from continuous driving scenes using deep neural networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 1, pp. 41–54, 2020.
- [31] P. Baldi, "Autoencoders, unsupervised learning, and deep architectures," in *ICML*, vol. 27, 2012, pp. 37–50.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [33] J. Hu, X. Guo, J. Chen, G. Liang, F. Deng, and T. L. Lam, "A two-stage unsupervised approach for low light image enhancement," *IEEE Robotics Autom. Lett.*, vol. 6, no. 4, pp. 8363–8370, 2021.
- [34] Y. Yang, B. Li, P. Li, and Q. Liu, "A two-stage clustering based 3d visual saliency model for dynamic scenarios," *IEEE Trans. Multim.*, vol. 21, no. 4, pp. 809–820, 2019.
- [35] A. Singh, H. Liu, G. Zhou, A. Yu, N. Rhinehart, and S. Levine, "Parrot: Data-driven behavioral priors for reinforcement learning," in *ICLR*, 2021.
- [36] C. Tan, F. Sun, B. Fang, T. Kong, and W. Zhang, "Autoencoder-based transfer learning in brain-computer interface for rehabilitation robot," *Int. J. Adv. Robot. Syst.*, vol. 16, no. 2, 2019.
- [37] D. Zadok, D. McDuff, and A. Kapoor, "Modeling affect-based intrinsic rewards for exploration and learning," in *ICRA*, 2021, pp. 13 078–13 084.
- [38] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.
- [39] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [40] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.
- [41] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *ICCV*, 2015, pp. 4489–4497.
- [42] A. M. Hari and L. Rajan, "Advanced materials and technologies for touch sensing in prosthetic limbs," *IEEE Trans. Nanobiosci.*, 2021.
- [43] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*, vol. 9351, 2015, pp. 234–241.