

Surgical Procedure Understanding, Evaluation, and Interpretation: A Dictionary Factorization Approach

Abed Soleymani, Xingyu Li, and Mahdi Tavakoli

Abstract—In this study, we present a novel machine learning-based technique to help surgical mentors assess surgical motion trajectories and corresponding surgical skills levels in surgical training programs. The proposed method is a variation of sparse coding and dictionary learning that is straightforward to optimize and produces approximate trajectory decomposition for structured tasks. Our approach is superior to existing stochastic or deep learning-based methods in terms of transparency of the model and interpretability of the results. We introduce a dual-sparse coding algorithm which encourages the elimination of redundant and unnecessary atoms and targets to reach the most informative dictionary, representing the most important temporal variations within a given surgical trajectory. Since surgical tool trajectories are time series signals, we further incorporate the idea of floating atoms along the temporal axis in trajectory analysis, which improves the model’s accuracy and prevents information loss in downstream tasks. Using JIGSAWS data set, we present preliminary results showing the feasibility of the proposed method for clustering and interpreting surgical trajectories in terms of user’s skills-related behaviors.

Index Terms—Machine Learning, Dictionary Learning, Sparse Coding, Surgical Trajectory Decomposition, Surgical Skills Assessment.

I. INTRODUCTION

Modern surgical robots are capable of measuring and recording surgical activities (e.g., kinematics data, video recordings, eye-gaze data), which makes them a perfect platform for incorporating data-driven solutions for procedural understanding and user skills assessment applications. There are several work that perform the autonomous robotic surgical skills evaluation using deep learning (DL) [1]–[5]. Even though these studies have reported relatively low misclassification rates, they are black-box models in which the decision-making procedure is not transparent or understandable in human terms. As a result, there is neither intuitive nor explainable feedback to the user about his/her surgical performance and contributing factors to the predicted model outcome (in this paper, *explainability* and *intuitiveness* are referred to as the extent to which the internal mechanism or outcome of a model can be explained and are intuitive in human terms, respectively). Moreover, the high capacity of the mentioned models (i.e., the ability of the model for accommodating input data variations which mainly depends on the number of learnable parameters), especially DL models, demands large training sets to avoid overfitting. Since in the field of robotic surgery, clean and reliable data sets are very small in size, such

models usually tend to overfit and fail to generate reliable models that perform well in novel situations (e.g., aborting and restarting a task, unwanted mistakes caused by poor depth estimation, etc.) [6].

To provide users with more elaborated targeted feedback about their skills level and surgical performance (e.g., in which part of the task the surgeon should improve his/her skills), one can break up surgical trajectories into pre-defined segments called *surgemes* [7] and apply surgical skills assessment methods at the sub-task level. In this paradigm, rather than having a global performance score, we will analyze the surgical workflow that results in high-resolution feedback regarding different aspects and parts of the whole executive task. There is a rich body of literature trying to perform fine-grained analysis of surgical activities in an automatic way using DL [8]–[13], reinforcement learning (RL) [14], and Hidden Markov Models (HMMs) [15]. These approaches not only have the same problems raised from being a black-box model but also suffer from over-segmentation (i.e. predicting numerous insignificant action boundaries) and low accuracy rate that prevents them to make certain predictions, especially for rare or unseen events (e.g., mid-task failures) [7]. We infer that the over-segmentation problem arises from the fact that these models pay too much attention to the data local variations (i.e., microscale details), rather than global context. Moreover, to evaluate the segmentation accuracy of these methods, they heavily rely on manually made gesture annotations, which would be very laborious, time consuming, and prone to inter-annotator variation. More importantly, the mentioned approaches break up trajectories into segments without offering any explicit interpretation about the behaviors and the dexterity of the user in the sub-task level.

Statistical machine learning (ML) techniques, on the other hand, usually perform better than DL models on small training data sets. An intuitive and human-inspired ML approach can provide us with a more transparent and explainable solution for data-scarcity surgical decomposition task [16]. One intuition about *structured* tasks such as suturing trials, human walking cycles, parallel parking, etc., is that their main variations can be decomposed into finite components namely *general trends* and *seasonal patterns*. This concept is illustrated in Fig. 1 in which \mathbf{d}_1 is the average general trend and other components are the averaged dominant seasonal patterns for all 6 trials within the training data set. In this kind of decomposition, each component contains specific and important intra-trial temporal variations. Moreover, taking Fig. 1 as an example, each trial in training or test data sets can be individually reconstructed via a linear combination of components, i.e., $\text{Trial}_k = \sum_{i=1}^3 c_{ki} \mathbf{d}_i$ where $c_{ki} \in \mathbb{R}$ is the weight indicating the contribution strength of component \mathbf{d}_i in the reconstruction of k^{th} trial. If \mathbf{d}_i are trained on expert trajectories, gains c_{ki} for a given test trial can be used to determine the fidelity of the participant to the averaged *ideal* trend

[†] All authors are with the Electrical and Computer Engineering Department, University of Alberta, Edmonton, Alberta, Canada. {zsoleymani, xingyu, mahdi.tavakoli}@ualberta.ca

[‡] This research was supported by the Canada Foundation for Innovation (CFI), UAlberta Huawei-ECE Research Initiative (HERI), the Government of Alberta, the Natural Sciences and Engineering Research Council (NSERC) of Canada, the Canadian Institutes of Health Research (CIHR), and the Alberta Economic Development, Trade and Tourism Ministry’s grant to Centre for Autonomous Systems in Strengthening Future Communities.

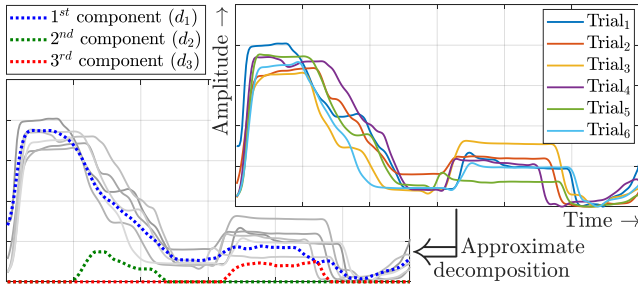


Figure 1: Approximate trajectory decomposition in structured tasks.

and seasonal patterns and convey important information regarding the execution quality and dexterity of the user. That is, each trajectory can be characterized by its coefficient vector $\mathbf{c}_k = [c_{k1}, c_{k2}, c_{k3}]^\top$ in the *embedding space* in which hidden behaviors of the user will be revealed.

Inspired by these intuitions, *dictionary learning* and *sparse coding* [17] can be modified to be applied on structured time series for meaningful, interpretable, and human-inspired trajectory decomposition task. Generated components (i.e., *dictionary atoms*) and their contribution in reconstruction of each individual trajectory (i.e., generated *code* matrix) disclose important information for several downstream tasks such as skills assessment, skills transfer, and anomaly detection.

In this research, we will introduce *dual-sparse* dictionary learning approach for the approximate trajectory decomposition of structured tasks in retrospective studies. We will also introduce our dual-sparse dictionary learning algorithm with a novel absolute mutual incoherence metric μ^+ . The idea of *floating atoms* will be incorporated in the proposed algorithm to accommodate trajectory structures with temporal shift. This preserves the relative temporal structure of the underlying events and prevents information loss while mapping the data to lower dimensional space (or embedding space, e.g., three dimensional space which is perceptible for humans and suitable for data visualization purposes) and makes our embedding representation more meaningful and realistic. Finally, we will evaluate our method on basic structured robotic surgery trajectories in JIGSAWS data set for surgical skills assessment and anomaly detection tasks. We believe that our novel approach has a potentially high impact in robotic surgery, where demands for enhanced safety and explainability are extremely strong.

The paper is organized as follows: In Section II, basic concepts and motivations that lead us to our contributions will be discussed. In Section III, algorithm implementation and the idea of floating atoms will be presented. In Section IV, fundamental features of our approach and its application on JIGSAWS data set will be investigated. In Section V, several discussions about the advantages and practical details of the proposed method will be presented. Concluding remarks are provided in Section VI.

II. PROBLEM STATEMENT

A. Preliminaries

Sparse coding is a method of representing data vectors as sparse linear combinations of a set of basis elements called *atoms*. It is as-

sumed that all atoms together which make *dictionary* matrix, capture main directions in the input space and have enough information for reconstructing input data. Dictionary learning algorithms try to develop a dictionary matrix that efficiently reconstructs each input data by a linear combination of the generated atoms. In this context, all generated coefficients of linear combinations are referred to as *code*.

A traditional dictionary learning framework for sparse representation optimizes the empirical loss function

$$\mathcal{L}(\mathbf{D}, \mathbf{C}) = \min_{\mathbf{D}, \{\mathbf{c}_i\}_{i=1}^n} \sum_{i=1}^n \left[\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 + \alpha \|\mathbf{c}_i\|_0 \right] \quad (1)$$

for the finite set of n data vectors $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$, aiming to find an optimal dictionary $\mathbf{D} := [\mathbf{d}_1, \dots, \mathbf{d}_p] \in \mathbb{R}^{d \times p}$ such that each data vector \mathbf{x}_i can be well-approximated by a linear combination of dictionary atoms $\{\mathbf{d}_j\}_{j=1}^p$. The term $\|\mathbf{c}_i\|_0$ in (1) which is the number of non-zero elements in vector \mathbf{c}_i , encourages to minimize the number of non-zero element in code vectors $\mathbf{c}_i \in \mathbb{R}^p$ of the code matrix $\mathbf{C} := [\mathbf{c}_1, \dots, \mathbf{c}_n] \in \mathbb{R}^{p \times n}$. The sparsity-promoting loss $\|\mathbf{c}_i\|_0$ encourages the coding algorithm to rely on a minimum number of the generated dictionary atoms for reconstruction which has the advantage of being simple and easy to decode, either by machine or human. The loss function (1) simply tries to make a balance between data reconstruction loss $\|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2$ and sparsity-promoting loss $\|\mathbf{c}_i\|_0$ with a regularization parameter α .

The l_0 sparsity loss $\|\mathbf{c}_i\|_0$ in (1) makes the optimization an NP-hard problem with sub-optimal solution [18]. Replacing $\|\mathbf{c}_i\|_0$ with its l_1 convex relaxation $\|\mathbf{c}_i\|_1$ yields optimal and sparse solutions for codes \mathbf{c}_i [19]. To prove why l_1 regularization term $\|\mathbf{c}_i\|_1$ enforces elements of vector \mathbf{c}_i to be zero (i.e., to be sparse rather than small), we encourage the reader to see [20].

Due to the bi-linearity between the dictionary \mathbf{D} and codes \mathbf{c}_i , (1) is still non-convex and cannot be jointly optimized with respect to dictionary \mathbf{D} and code matrix \mathbf{C} [21]. Since (1) is convex with respect to variables \mathbf{D} or \mathbf{C} when the other one is fixed, a solution to this problem is to alternate the optimization procedure between the two variables, i.e., minimizing the loss function with respect to one parameter while keeping the other one fixed. To prevent dictionary atoms to get arbitrary large in the optimization process, we should normalize each atom (i.e., $\mathbf{d}_i \leftarrow \mathbf{d}_i / \|\mathbf{d}_i\|_2, \forall i$) after each dictionary update.

Conventional dictionary learning approach usually generates more atoms than the number of data samples and creates over-complete dictionaries (e.g., $p \approx 4n$ in [22]) for the sake of reconstruction accuracy. Besides the high computational cost and memory demand of the over-complete dictionary approaches, the redundancy of the generated dictionary does not necessarily enhance the optimality and performance of the final solution. One major problem with over-complete dictionaries is the high correlation between generated atoms, which degrades the *mutual incoherence* metric defined as

$$\mu(\mathbf{D}) = \max_{i \neq j} \frac{|\mathbf{d}_i^\top \mathbf{d}_j|}{\|\mathbf{d}_i\|_2 \|\mathbf{d}_j\|_2}. \quad (2)$$

It is empirically observed that highly correlated atoms in over-complete dictionaries or in general, high values for $\mu(\mathbf{D})$ (in Section III-A we will calculate an upper bound for μ) make the sparse

coding stage slow, computationally demanding, and non-optimal [23]. Tackling this problem, [24] proposes orthogonal dictionary learning method (we name it $\mathcal{L}^{\text{orth}}(\mathbf{D}, \mathbf{C})$) that is the constrained version of (1) subject to $\mathbf{D}^\top \mathbf{D} = \mathbf{I}_{p \times p}$ where \mathbf{I} indicates identity matrix. [24] argues that this approach yields $\mu(\mathbf{D}) = 0$, faster convergence, and comparable results relative to other sophisticated over-complete dictionary learning methods such as K -SVD [25].

B. Motivation

Prior work in surgical skills evaluation has shown that a lay observer is able to discover and rate the skillful behavior of a surgeon just by looking at his/her pre-recorded translational and/or rotational hand movement patterns with accuracy comparable to an expert surgeon [26]. A possible explanation about this interesting result is that the lay observer does not care about extreme details and microscale translations/rotations within the surgical trajectory. He/she just pays attention to the most informative temporal features within the executive task, i.e., (1) general trend, (2) seasonal patterns, and (3) unwanted random actions. Although the cognitive procedure of decision making based on these three factors is unknown, we are inspired to investigate how much the fidelity to the general trend, correct execution of each seasonal pattern, and minimum occurrence of incidental motions play a crucial role in rating the performance of the user in surgical tasks.

The central idea of this work is to generate an intuitive and universally understandable representation of surgical trajectories based on building blocks of surgery that meaningfully describes the procedural flow and highlights hidden information of a surgical task. The core intuition is aligned with the motivation of representing data as a sparse linear combination of specific atoms in sparse coding problems. In this context, each dictionary atom can be a representative of the task general trend or a seasonal pattern (i.e., surgeme) that encapsulates one important local variation of the trajectory.

Unlike prior dictionary learning algorithms that rely on generating a lot of atoms to take care of details and microscale variations of data, in this work, we develop an algorithm that selectively removes unnecessary atoms and focuses on preserving important variations within the input data to achieve an understandable and interpretable trajectory decomposition. Apart from neglecting unnecessary details and the small number of atoms, another feature that makes our representation meaningful and easy to interpret is the minimum amount of overlap between two arbitrary atoms. This is meaningful because we aim to assign each non-overlapping seasonal pattern of the structured trajectory to one dictionary atom (e.g., 2nd and 3rd components in Fig. 1). Moreover, a small overlap between dictionary atoms indicates that the action executed in a particular timestamp can be purely attributed to one dominant atom and it simplifies interpretations about the quality of the task done in that timestamp. Factorizing minimally-overlapping atoms can be thought of performing approximate decomposition for a particular structured task.

Although enforcing atoms to have zero overlap with each other during the learning process makes them more explainable, it degrades the signal reconstruction quality with poor non-smooth results and the interpretability of generated codes. Moreover, as we will show later, intensively reducing the total number of active atoms increases their overlap. This is because the algorithm

tries to allocate all variations of the trajectory among currently existing atoms to minimize the reconstruction loss. There is an inter-dependency between the number of atoms, their overlap, and reconstruction loss to generate informative codes for a given set of trajectories. We call the smallest number of minimally-overlapping atoms that gives us a relatively good reconstruction the *intrinsic dimensionality of embedding space* (δ).

C. Data Set

All analysis in this work are done based on the standard JIGSAWS data set [27] collected from surgical activities of eight surgeons in three different levels of expertise (i.e., novice, intermediate, and expert) performing suturing (SU), knot-tying (KT), and needle-passing (NP) tasks on the *da Vinci* Surgical System. JIGSAWS contains three Cartesian motions along x , y , and z axes as well as 9 elements of rotational matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ for both hands of the user and also for both patient-side robotic arms. Note that, all 9 elements of rotational matrix \mathbf{R} can be expressed as 3 angles *roll* (Φ), *pitch* (Θ), and *yaw* (Ψ) as follows

$$\Phi = \text{atan2}\left(\frac{r_{21}}{r_{11}}\right), \Theta = \text{atan2}\left(\frac{-r_{31}}{\sqrt{r_{32}^2 + r_{33}^2}}\right), \Psi = \text{atan2}\left(\frac{r_{32}}{r_{33}}\right)$$

where r_{ij} is the element in the i^{th} row and j^{th} column of \mathbf{R} .

III. METHODOLOGY

A. Dictionary Factorization

We will show that having the minimum overlap between two atoms is a much stricter condition compared to the orthogonality condition (i.e., minimum correlation between atoms or small value for $\mu(\mathbf{D})$) presented in [24]. To have a measure of overlap between two arbitrary atoms, we introduce the *absolute mutual incoherence* metric defined as

$$\mu^+(\mathbf{D}) = \max_{i \neq j} \frac{|\mathbf{d}_i^\top \mathbf{d}_j|}{\|\mathbf{d}_i\|_2 \|\mathbf{d}_j\|_2}. \quad (3)$$

where $|\mathbf{d}_i|$ denotes the element-wise absolute value of the vector \mathbf{d}_i . Since $|\mathbf{d}_i^\top \mathbf{d}_j| \geq |\mathbf{d}_i^\top \mathbf{d}_j|$ holds for $\forall \mathbf{d}_i, \mathbf{d}_j \in \mathbb{R}^d$, from (2) and (3) it can be concluded that $\mu^+(\mathbf{D})$ is an upper bound for $\mu(\mathbf{D})$ (i.e., $\mu(\mathbf{D}) \leq \mu^+(\mathbf{D})$) and minimizing $\mu^+(\mathbf{D})$ shrinks $\mu(\mathbf{D})$. However, minimizing $\mu(\mathbf{D})$ does not necessarily yield reduced $\mu^+(\mathbf{D})$.

Property 1: $0 \leq \mu(\mathbf{D}) \leq \mu^+(\mathbf{D}) \leq 1$. □

Proof. See Appendix A. ■

We argue that minimizing the total overlap between atoms promotes the reduction of $\mu^+(\mathbf{D})$. To formulate overlapping between atoms, we rewrite the dictionary matrix based on its rows, $\mathbf{D} := [\hat{\mathbf{d}}_1^\top, \dots, \hat{\mathbf{d}}_d^\top]^\top$ where $\hat{\mathbf{d}}_k = [\mathbf{d}_1(k), \dots, \mathbf{d}_p(k)]$ is the k^{th} row of the matrix \mathbf{D} and $\mathbf{d}_i(k)$ is the k^{th} element of atom \mathbf{d}_i . Note, applying l_1 regularization on $\hat{\mathbf{d}}_k$ (i.e., minimizing $\|\hat{\mathbf{d}}_k\|_1 = \sum_{i=1}^p |\mathbf{d}_i(k)|$) in a quadratic objective function (e.g., $\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2$) enforces elements of $\hat{\mathbf{d}}_k$ (i.e., $\mathbf{d}_i(k)$) to be zero which eventually reduces the total overlap between all atoms at timestamp k . This means setting $|\mathbf{d}_i(k)| |\mathbf{d}_j(k)|$ to zero for as many as possible $1 \leq i, j \leq p$, $i \neq j$, which ultimately minimizes $\mu^+(\mathbf{D})$. Note that since $\mu^+(\mathbf{D})$ normalizes each atom, reducing $|\mathbf{d}_i(k)| |\mathbf{d}_j(k)|$ is not enough; it is ideal for it to be zero. It means we

have to have $\mathbf{d}_i(k)$ equal to zeros for a lot of all possible i for a given k (i.e., a sparse $\hat{\mathbf{d}}_k$ for all $1 \leq k \leq d$), which conceptually reduces the overlap between atoms at timestamp k . As a result, $\mu^+(\mathbf{D})$ is a good measure of overlap between dictionary atoms which $\mu^+(\mathbf{D})=1$ for a fully-overlapped dictionary and $\mu^+(\mathbf{D})=0$ for non-overlapping dictionary. Later, we will investigate the relationship between $\mu^+(\mathbf{D})$ and the total overlap between active dictionary atoms.

1) *Objective Function Definition:* Now, we define a new cost function $\mathcal{L}^+(\mathbf{D}, \mathbf{C})$ that tries to jointly generate sparse codes and dictionaries with minimally-overlapping atoms to reconstruct the input data \mathbf{X}

$$\mathcal{L}^+(\mathbf{D}, \mathbf{C}) = \min_{\mathbf{D}, \{\mathbf{c}_i\}_{i=1}^n} \left\{ \sum_{i=1}^n \left[\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 + \alpha \|\mathbf{c}_i\|_1 \right] + \sum_{k=1}^d \beta \|\hat{\mathbf{d}}_k\|_1 \right\} \quad (4)$$

where β is the regularization factor that determines the relative importance of total overlap compared to the sparsity of generated codes and lumped reconstruction error for all data vectors in \mathbf{X} .

2) *Algorithm Implementation:* As explained in Section II-A, for the sake of convexity we need to alternatively optimize loss function (4) with respect to one parameter \mathbf{D} or \mathbf{C} while keeping the other one fixed. As a result, we have two steps for minimizing our cost function: \mathbf{C} -step and \mathbf{D} -step. In \mathbf{C} -step we fix matrix \mathbf{D} (coming from initialization or the previous \mathbf{D} -step) and do

$$\mathcal{L}_{\mathbf{C}}^+(\mathbf{D}, \mathbf{C}) = \min_{\{\mathbf{c}_i\}_{i=1}^n} \sum_{i=1}^n \left[\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 + \alpha \|\mathbf{c}_i\|_1 \right]. \quad (5)$$

This optimization is equivalent to the *lasso* problem [28] for which many fast algorithms exist. We assume that the function $\text{sparseCode}(\mathbf{X}, \mathbf{D}, \alpha)$ gets the dictionary \mathbf{D} , data vectors \mathbf{X} , and regularization parameter α and returns code matrix \mathbf{C} .

In \mathbf{D} -step, we fix matrix \mathbf{C} and do

$$\mathcal{L}_{\mathbf{D}}^+(\mathbf{D}, \mathbf{C}) = \min_{\mathbf{D}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 + \sum_{k=1}^d \beta \|\hat{\mathbf{d}}_k\|_1. \quad (6)$$

Due to our l_1 condition on the rows of dictionary matrix \mathbf{D} , (6) is no longer a conventional dictionary update problem with closed-form solution via coordinate descent approach. Moreover, we have two summations that make the solution formulation hard and long.

We define the reconstruction error vector for each data sample \mathbf{x}_i as $\boldsymbol{\varepsilon}_i := \mathbf{x}_i - \mathbf{D}\mathbf{c}_i \in \mathbb{R}^d$ and reconstruction error matrix as $\mathbf{E} := [\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n] = \mathbf{X} - \mathbf{D}\mathbf{C} \in \mathbb{R}^{d \times n}$. According to the norm-2 definition, $\|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 = \|\boldsymbol{\varepsilon}_i\|_2^2 = \sum_{j=1}^d e_{ji}^2$ where e_{ji} is the element in the j^{th} row and i^{th} column of matrix \mathbf{E} . As a result

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 = \sum_{i=1}^n \left(\sum_{j=1}^d e_{ji}^2 \right) = \sum_{k=1}^d \sum_{r=1}^n \hat{e}_{rk}^2 \quad (7)$$

where \hat{e}_{rt} is the element in the r^{th} row and k^{th} column of matrix \mathbf{E}^{\top} . (7) shows the trivial fact that the sum of squares of all elements in matrix \mathbf{E} is equal to that of \mathbf{E}^{\top} . Due to the fact that

Algorithm 1: Dual-sparse coding algorithm.

Input: \mathbf{X} , α , and β

Result: Dictionary \mathbf{D} with

minimally-overlapping atoms and code matrix \mathbf{C}

Call function $\text{sparseCode}(\cdot, \cdot, \cdot)$

Initialize dictionary \mathbf{D}

while not converged do

 # \mathbf{C} -step

$\mathbf{C} \leftarrow \text{sparseCode}(\mathbf{X}, \mathbf{D}, \alpha)$

 # \mathbf{D} -step

$\mathbf{D} \leftarrow \left[\text{sparseCode}(\mathbf{X}^{\top}, \mathbf{C}^{\top}, \beta) \right]^{\top}$

end

#normalizing atoms of dictionary \mathbf{D}

for $i \leftarrow 1$ **to** p **do**

if $\|\mathbf{d}_i\|_2 \neq 0$ **then**

$\mathbf{d}_i \leftarrow \mathbf{d}_i / \|\mathbf{d}_i\|_2$

else

 #removing nullified atoms

 Remove \mathbf{d}_i from \mathbf{D}

end

end

#generating finalized code matrix

$\mathbf{C} \leftarrow \text{sparseCode}(\mathbf{X}, \mathbf{D}, \alpha)$

return \mathbf{D} and \mathbf{C}

$\mathbf{E}^{\top} = (\mathbf{X} - \mathbf{D}\mathbf{C})^{\top} = \mathbf{X}^{\top} - \mathbf{C}^{\top}\mathbf{D}^{\top}$ we have

$$\sum_{k=1}^d \left(\sum_{r=1}^n \hat{e}_{rk}^2 \right) = \sum_{k=1}^d \|\hat{\mathbf{x}}_k - \mathbf{C}^{\top}\hat{\mathbf{d}}_k\|_2^2 \quad (8)$$

where $\hat{\mathbf{d}}_k$ is the k^{th} column of \mathbf{D}^{\top} (or as previously mentioned the k^{th} row of \mathbf{D}) and $\hat{\mathbf{x}}_k$ is the k^{th} column of \mathbf{X}^{\top} . Combining (7) and (8) results

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 = \sum_{k=1}^d \|\hat{\mathbf{x}}_k - \mathbf{C}^{\top}\hat{\mathbf{d}}_k\|_2^2. \quad (9)$$

Substituting (9) into (6) yields a more tractable cost function for our specific dictionary update procedure

$$\mathcal{L}_{\mathbf{D}^{\top}}^+(\mathbf{C}^{\top}, \mathbf{D}^{\top}) = \min_{\mathbf{D}^{\top}} \sum_{k=1}^d \left[\frac{1}{2} \|\hat{\mathbf{x}}_k - \mathbf{C}^{\top}\hat{\mathbf{d}}_k\|_2^2 + \beta \|\hat{\mathbf{d}}_k\|_1 \right]. \quad (10)$$

Interesting point about (10) is that it tries to minimize exactly the same problem described in (5). In fact, our specialized dictionary learning problem for a given code matrix \mathbf{C} and data set \mathbf{X} became an sparse coding problem for the dictionary \mathbf{C}^{\top} and data set \mathbf{X}^{\top} with regulation parameter β . The algorithm for solving the dual-sparse coding problem is described in Algorithm 1.

The first point regarding Algorithm 1 is that unlike the traditional dictionary learning methods, in our method norms of atoms do not arbitrary get large during the optimization process since the size of each atom is penalized by minimizing $\sum_{k=1}^d \|\hat{\mathbf{d}}_k\|_1$ in (4). As a result, we do not need to normalize dictionary atoms \mathbf{d}_i after each update. Second, we start the algorithm with initialized dictionary matrix with the number of atoms higher than the intrinsic

dimensionality of embedding space of that particular task to give the algorithm the chance of deleting unnecessary atoms in its own way. In the end, the final code matrix C will be generated based on the normalized final D .

3) *Algorithm Convergence*: In each iteration, Algorithm 1 aims to reduce the total reconstruction loss, code sparsity loss, and dictionary atoms overlapping loss by minimizing convex cost functions (5) and (10) in C and D steps, respectively. According to [29], [30], if C^* and D^* are optimal solutions for the loss function (4) and the Algorithm 1 starts from C_0 and D_0 , for each iteration $k \geq 0$ we have

$$\mathbb{E}[\mathcal{L}^+(\mathbf{D}_{k+1}, \mathbf{C}_{k+1})] - \mathcal{L}^+(\mathbf{D}^*, \mathbf{C}^*) \leq \xi(\mathcal{L}^+(\mathbf{D}_k, \mathbf{C}_k) - \mathcal{L}^+(\mathbf{D}^*, \mathbf{C}^*)) \quad (11)$$

where $0 \leq \xi < 1$ is a constant derived from the properties of loss function (4) (see Appendix B) and $\mathbb{E}[\mathcal{L}^+(\mathbf{D}_{k+1}, \mathbf{C}_{k+1})]$ is the expected value of loss $\mathcal{L}^+(\mathbf{D}, \mathbf{C})$ in $(k+1)^{\text{th}}$ iteration. (11) implies that in each iteration, statistically, it is guaranteed that the expected value of the loss function (4) approaches to its minimum value $\mathcal{L}^+(\mathbf{D}^*, \mathbf{C}^*)$ through the coordinate descent Algorithm 1.

B. Dictionary Temporal Fine-tuning

Some atoms that are generated by Algorithm 1 are similar to islands with sharp onsets, short duration, and no overlap with each other (e.g., green and black atoms in the left column of Fig. 2). These atoms are the result of averaging a significant variation (i.e., surgeme) within trajectories of the training set that the algorithm perceived they are important in reconstructing the original trajectory.

An issue that may arise with these specific atoms is that they can be arbitrarily aligned with respect to the structure of a given trajectory within the test set (i.e., they can appear at different temporal positions with some phase shifts within a given executive task). In other words, although these atoms are presented in the optimal place with regards to the trajectories of the training set, their temporal position might not be optimal for a particular trajectory (either from the training set or test set). This phenomenon can neglect important variations within the original trajectory during the reconstruction and coding stage and as a result, distorts the values of the generated codes for those atoms. When the high-dimensional data (e.g., trajectory data vector) is mapped to lower-dimensional embedding space, unwanted changes in code vectors can be treated as information loss. Since the hidden behaviors of each user are mapped to the embedding space, any sort of information loss deteriorates the accuracy of any interpretation and evaluation of that trajectory.

One possible solution to this problem is to perform a post-processing step in the training stage and shift the atoms slightly to the left and right and investigate where this atom fits best (this is why we call these atoms *floating* atoms). This type of modeling due to the extra degree-of-freedom on generating atoms benefits from the advantages of over-complete dictionaries in terms of having lower reconstruction and information losses and at the same time, due to the small number of minimally-overlapping atoms avoids disadvantages of over-complete dictionary learning methods such as the high correlation between atoms discussed earlier in Section II-A. The bounded shifting to left and right is meant to secure the floating atom for intended variation within the trajectory, and to prevent it from

Algorithm 2: Meta-algorithm for calculating optimal shift φ_i for floating atoms d_i .

Input: $\mathbf{x}_j, \mathbf{D}, \varphi_{\max}$, and γ
Result: φ_i
 Call function sparseCode
 $\varphi_i \leftarrow -\varphi_{\max}$ #initializing optimal shift
 $C_{\max} \leftarrow -\infty$ #initializing optimal cost
for $\varphi \leftarrow -\varphi_{\max}$ **to** φ_{\max} **do**
 $\mathbf{c}_{i,\varphi} = \text{sparseCode}(\mathbf{x}_j, \mathbf{D}_\varphi, \alpha)$
 $\mathcal{C} = \mathbf{x}_j \star \mathbf{d}_{i,\varphi} - \gamma \|\mathbf{x}_j - \mathbf{D}_\varphi \mathbf{c}_{i,\varphi}\|_2$
 if $\mathcal{C} > C_{\max}$ **then**
 $C_{\max} \leftarrow \mathcal{C}$
 $\varphi_i \leftarrow \varphi$
end
end
return φ_i

mixing with other nearby atoms. Later we will explain one heuristic for calculating the maximum amount of shifting (i.e., φ_{\max}).

Without loss of generality, assume that first f atoms of p atoms meet all conditions of floating atoms (i.e., island-shaped, sharp onset, low duration, and no overlap with other atoms except general trend atom). A good metric to find the optimal temporal shift for a floating atom (i.e., φ_i for floating atom d_i where $1 \leq i \leq f$) is cross-correlation. Higher cross-correlation between trajectory \mathbf{x}_j and floating atom d_i with specific time shift φ_i (we define it d_{i,φ_i}) means that the trajectory needs the atom to be presented in that shifted position for the better reconstruction and less information loss. If the amount of cross-correlation is the same for several shifts within the domain of φ_i (i.e., $[-\varphi_{\max}, \varphi_{\max}]$), low reconstruction loss for \mathbf{x}_j based on new dictionary \mathbf{D}_{φ_i} (i.e., dictionary \mathbf{D} with $d_i \leftarrow d_{i,\varphi_i}$ modification) and new generated sparse code \mathbf{c}_{i,φ_i} will finalize the optimal value for shifting. The optimal shift φ_i for floating atoms d_i given trajectory \mathbf{x}_j will be calculated by solving this optimization problem

$$\varphi_i = \underset{\varphi}{\operatorname{argmax}}(\mathbf{x}_j \star \mathbf{d}_{i,\varphi} - \gamma \|\mathbf{x}_j - \mathbf{D}_\varphi \mathbf{c}_{i,\varphi}\|_2) \quad (12)$$

where \star is cross-correlation operand, γ is the regularization factor, and $\varphi \in [-\varphi_{\max}, \varphi_{\max}]$. The meta-algorithm for calculating the optimal shifting for floating atoms that can be applied to all floating atoms of a learned dictionary is presented in Algorithm 2. After applying Algorithm 2 on all floating atoms d_i with respect to the test trajectory \mathbf{x}_j and finding φ_i for all possible $1 \leq i \leq f$, we will update dictionary \mathbf{D} to \mathbf{D}^\dagger . From now on, the generated code matrix $\mathbf{C}_j^\dagger = \text{sparseCode}(\mathbf{x}_j, \mathbf{D}^\dagger, \alpha)$ will be used in downstream applications discussed in the future sections.

Each row of Fig. 2 shows one example of how floating atoms can be incorporated to reduce the information loss and enhance the accuracy of the time series mapping for a given structured task. Dictionary atoms are plotted in the left column of Fig. 2 in which d_1^f is the general trend and green and black diagrams are two atoms capturing seasonal patterns of the trajectory that satisfy three conditions of being floating atoms. Fig. 2(b) and Fig. 2(e) show the reconstruction quality and relative temporal positioning of each

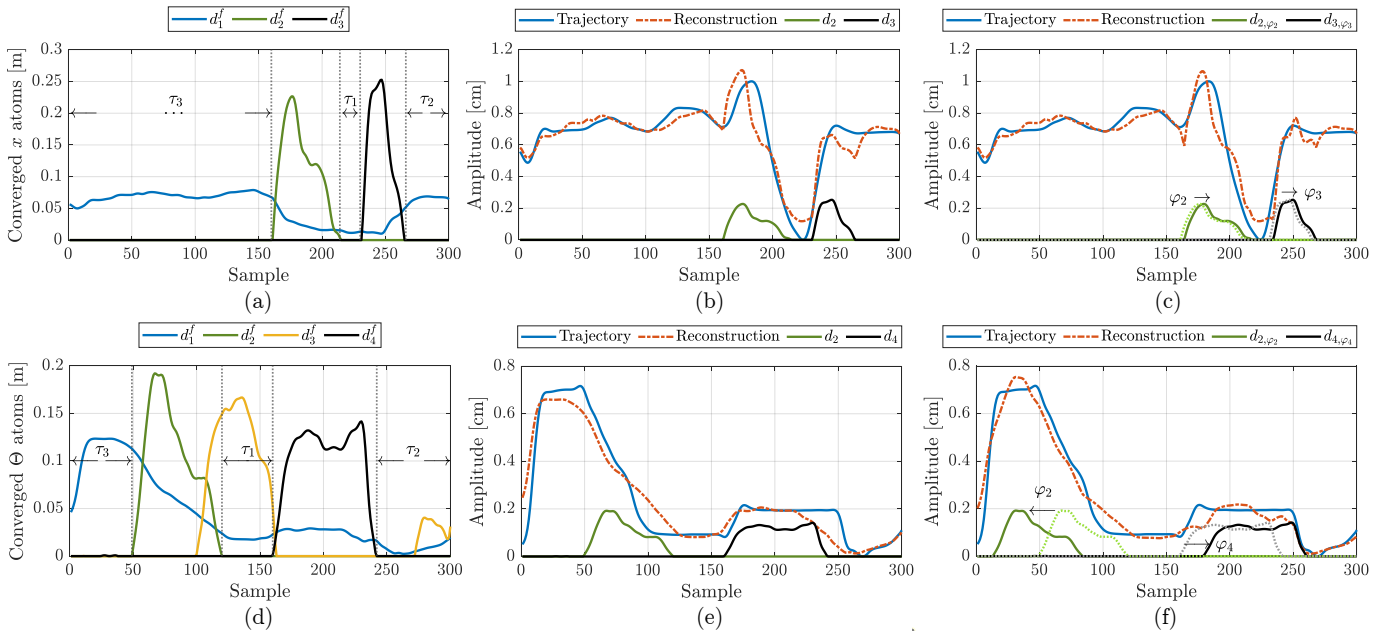


Figure 2: The effect of temporal fine-tuning of floating atoms of a given structured task on the reconstruction and information losses, maximum possible temporal shifts (τ_i), and optimal temporal shifts (φ_i). (a) and (d) show dictionary atoms that blue lines are general trend and green and black lines are two floating atoms for that specific trajectory. (b) and (e) show bad reconstruction and information losses due to the ill temporal positioning of the floating atoms. (c) and (f) illustrate the optimal positioning of floating atoms with regards to a given trajectory.

floating atom with respect to the original trajectories. As it is clear in Fig. 2(b), it is better if d_2^f and d_3^f slightly shift to the right to better capture the temporal variation that they have to represent. As shown in Fig. 2(c), optimal shift values are $\varphi_2 = \varphi_3 = 2$ samples. This change will result in 18.7% improvement of reconstruction loss and higher value for the codes assigned to these floating atoms (i.e., c_2 and c_3) that make them more accurate and meaningful in terms of reflecting skills and hidden behaviors of the user. The effect of floating atoms on Fig. 2(f) is even more noticeable. By shifting d_2^f for $\varphi_2 = 36$ to the left and d_4^f for $\varphi_4 = 19$ to the right (see Fig. 2(f)), we will achieve 27.4% improvement in reconstruction loss and considerable changes in the value of codes c_2 and c_4 . For instance, the temporal position of d_2^f with respect to the given trajectory in Fig. 2(e) was too bad that the sparse coding algorithm decided to neglect this atom in reconstruction and set c_2 to zero. After fine-tuning d_2^f in 2(f), the generated code for d_2^f became $c_2 = 1.11$ which sounds more accurate and realistic in human terms.

An important objective during shifting floating atoms is to avoid merging these atoms (i.e., they should not overlap after being placed in their optimal temporal positions). For instance, in Fig. 2(a), we should have $\varphi_{\max} \leq \tau_1/2$ to avoid merging two floating atoms in their extreme shifted states. In Fig. 2(d), d_2^f and d_4^f should not move to right and left, respectively to avoid increasing their overlap with d_3^f . These objectives not only preserve the isolated nature of floating atoms but also reduce the computational cost of finding optimal values for all possible shifts. If we maintain the non-overlapping property of floating atoms while shifting them to find optimal φ_i , they will remain decoupled. As a result, we can independently find φ_i for each valid d_i instead of jointly optimizing

(12) for all possible floating atoms.

Another upper bound for φ_{\max} is the maximum amount of sliding for each floating atom until they reach left or right boundaries of the time-series (e.g., τ_2 and τ_3 for the black and green atoms respectively in the left column of Fig. 2). A possible heuristic for finding a good upper bound for φ_{\max} based on explanations above and examples shown in Fig. 2(a) and Fig. 2(d) are $\varphi_{\max} \leq \min\{\tau_1/2, \tau_2, \tau_3\}$ and $\varphi_{\max} \leq \min\{\tau_2, \tau_3\}$, respectively.

Finding an appropriate value for φ_{\max} is case-dependent and needs human investigations for getting acceptable results. Since φ_{\max} should be calculated for each floating atom of each component within all trajectories (e.g., x , y , z , Φ , Θ , and Ψ of left and right hands), the complexity of the dictionary temporal fine-tuning is proportional to the complexity of the task (i.e., having numerous sub-tasks) and the number of components in each trajectory.

IV. APPLICATION TO JIGSAWS DATA SET

Inspired by the fact that decomposing surgical trajectories into their main variations (i.e., surgemes) reflects skills better than methods based on execution time or total path length [31], we resampled all surgical trials within the JIGSAWS data set to 300 samples and rescaled them between 0 and 1. According to our investigations, no data variation is removed during the resampling/rescaling since we have no sudden motion in surgical tasks. We can investigate task execution time and total path length as other two factors for our further investigations.

A. Model Training

Results presented here and in the next section are based on atoms trained over surgical data of an expert user out of 8 subjects in

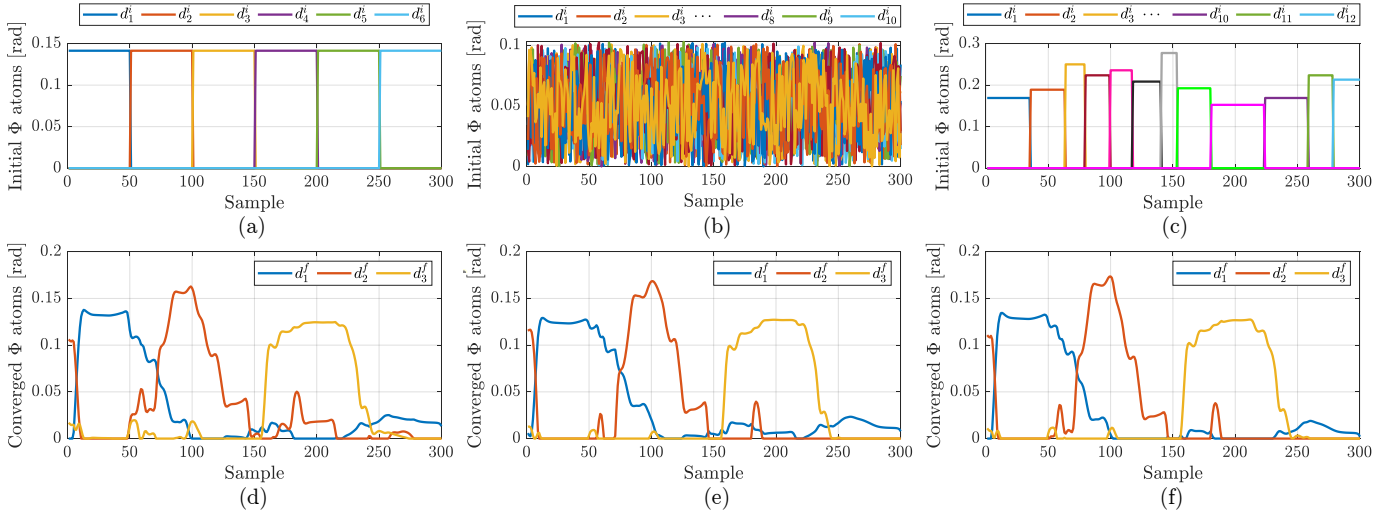


Figure 3: The effect of different initialization of dictionary matrix D in Algorithm 1 convergence. Each column shows the initialized atoms at the top and their converged final atoms with $\delta=3$ at the bottom. Different initializations converged to almost identical results (especially in the main variations) which indicates the consistency of the dual-sparse algorithm in finding meaningful atoms.

the JIGSAWS data set. Using expert data for training is due to the fact that global trends and seasonal patterns with minimum random movements within a particular task can be found in expert trajectories. Such a model can be used as a benchmark for other users to discover their hidden abnormal behaviors. It is worth mentioning that the processing time for different initializations and different trajectory components ranges from about 15 to 25 seconds.

B. Convergence Consistency

In a conventional dictionary learning approach the algorithm converges to very different dictionaries (i.e., non-similar local optima) for differently initialized atoms. Proper, robust, and informative initial point as an important factor for the success of the dictionary learning algorithm is an intensive field of research [32]. Due to our motivation, which is to meaningfully interpret an atom as a representative of a surgeme or sub-task of a surgical trajectory, randomly generated atoms are practically unusable even if they give us a good reconstruction.

In our setting, due to the dual-confined loss function described in (4), the small number of final atoms, and the structured nature of input time-series (that are coming from basic structured surgical tasks), Algorithm 1 tends to generate very similar atoms for different initial points for the dictionary. As it is clear in Fig. 3, different initializations with the different number of atoms p and different temporal lengths converged to almost identical atoms (especially in the main variations) with negligible difference in low amplitude parts. As we discussed in Appendix B, since (4) is *directional* component-wise Lipschitz continuous gradient, D^* and C^* in (11) are not necessarily global optima and hence, the convergence towards local optima is guaranteed in Algorithm 1. According to Fig. 3, the consistency of Algorithm 1 in converging to almost identical results for completely different initializations indicates that local optimum solutions are very close to the global optimum solution of (4).

The reasoning behind this consistency is that the cost function (4) *nullifies* unnecessary atoms (i.e., setting them to a vector of

zeros) while forming residual ones to capture main variations within the training trajectories that represent surgements and basic actions in a particular task. It is observed that Algorithm 1 tries to leave important variations of the currently nullified atom as inheritances to its neighboring active atoms to satisfy the reconstruction loss function. We argue that this behavior is the key element that makes our approach robust against the effect of dictionary initialization.

C. Reconstruction Quality

Perfect reconstruction (i.e., capturing almost all microscale details of input) is a major objective in the classic dictionary learning problem and also one of motivations for incorporating over-complete dictionaries. In our setting, perfect reconstruction requires extra atoms to represent unwanted random actions within structured trajectories which increases the risk of overfitting the training set and makes it hard to interpret the cause and effect of generated codes and explain them to a human. Inspired by this fact, the motivation of our approach is to sacrifice perfect reconstruction to achieve a small number of minimally-overlapping atoms that represent prevailed ongoing surgements within the trajectory. Following interpretations about the model reconstruction behavior in Fig. 5 indicate that this sacrifice is not in vain and helps to create an explainable approach for trajectory assessment.

Both trajectories in Fig. 5 and trajectory shown in Fig. 4(a) are reconstructed according to the atoms shown in Fig. 3(d) with hyperparameters $\alpha = 1$, $\beta = 1.5$, and $\gamma = 0.5$. The upper plot of Fig. 5 belongs to an intermediate user performing suturing trial with the code matrix of $c_{In} = [7.33, 4.89, 3.13]^T$ which is the same task performed by the expert user shown in Fig. 4(a) with the code matrix of $c_{Ex} = [5.89, -0.57, 7.26]^T$. One notable source of difference is the unusual behavior at the middle part of the intermediate trajectory that is the sign of happening mid-task failures and restarting while inserting the suture needle inside the phantom tissue. This common anomaly in surgical tasks is noticeable via bad reconstruction and will leave its trace in embedding space by distorting the code

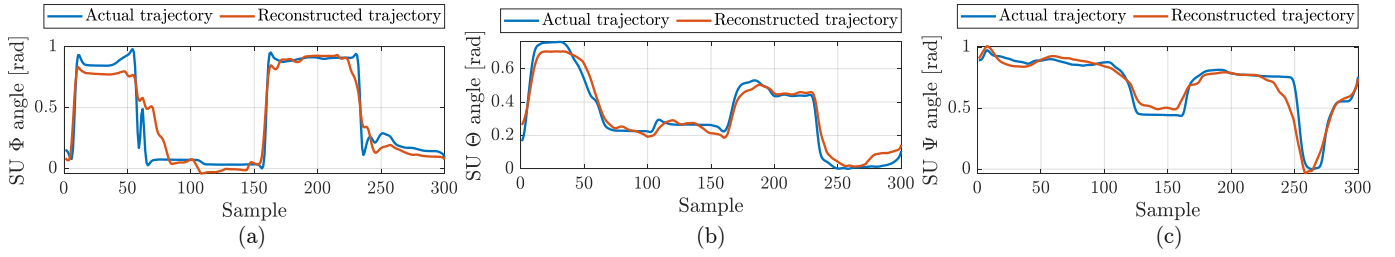


Figure 4: Reconstruction quality of the proposed method for rotational angles of the suturing task in JIGSAWS data set. We aim to preserve main variations of trajectories while neglecting unnecessary microscale details to create explainable atoms for the decomposition task.

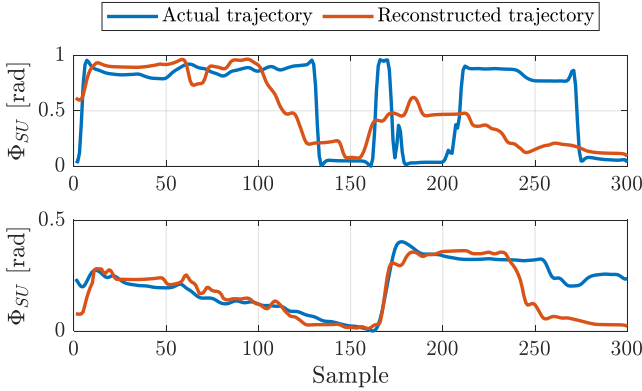


Figure 5: The effect of fault and lack of expertise on reconstruction quality of two sample trajectories in suturing task in JIGSAWS.

values of its neighbouring atoms. In this particular case, c_2 for instance, should be a low value for normal task execution (e.g., $c_2 = -0.57$ in c_{Ex}), but because of the anomaly, c_2 becomes higher than usual in c_{In} . This example indicates that although the effect of random actions is not explicitly a part of problem formulation, they implicitly leave their interpretable tracks in code space.

The lower plot of Fig. 5 belongs to a novice user with the code matrix of $c_{No} = [1.92, 0.59, 2.49]^T$ which is significantly dissimilar to the expert trial shown in Fig. 4(a) with the code matrix of c_{Ex} . For instance, one important source of the dissimilarity is the low values of roll rotation at the beginning of the trial that results in low value for code c_1 corresponding to d_1^f in Fig. 3(d). The low value of c_1 is the sign of a fundamental mistake in suturing task that will be elaborated in Section IV-D. For the sake of further clarification, a supplementary video is provided that explains the descriptions above along with the endoscopic videos of the trials plotted in Fig. 5.

D. Embedding Space Analysis

In a broader sense, the proposed dual-sparse dictionary learning approach maps the input trajectory into a lower-dimensional space (i.e., code or embedding space) that reveals the latent temporal structure of data that can be used for skills assessment, anomaly detection, and educational purposes. For the sake of clarification, embedding space of Φ angle for all suturing trials based on atoms in Fig. 3(d) trained over sample trajectories of the first expert Ex_1 is shown in Fig. 6(a) and Fig. 6(d). Similarly, embedding space of Θ angle for all suturing trials based on atoms in Fig. 2(d) trained

over sample trajectories of Ex_1 is shown in Fig. 6(b) and Fig. 6(e). Finally, embedding space of Ψ angle for all suturing trials based on dictionaries trained over sample trajectories of Ex_1 is shown in Fig. 6(c) and Fig. 6(f). As it is illustrated in Fig. 6, representations of different trials of each subject cluster near to each other in the embedding space. This behavior reflects their surgical *style* in the low dimensional space. Moreover, subjects with similar skills levels usually cluster near to each other. This can be useful for investigating the learning curve of a trainee while his/her latent representation approaches towards expert clusters after few sessions of training.

Another interesting point about plots in Fig. 6 is that, most non-expert subjects are pretty similar in deviating from the general trend d_1^f (i.e., $c_1 \approx 0$). This behavior results in compact near clusters along the c_1 axis for less-experienced users. However, since each non-expert user has his/her own way of deviating from the general trend this coherency does not imply consistency in performing the task. On the other hand, expert users have higher fidelity to the general trend by having large codes (e.g., $4 < c_1 < 8$ in Fig. 6(d)). However, this variation in c_1 cannot be attributed to the lack of consistency in performing the task. This is because d_1^f is normalized to one and has a relatively small maximum value compared to that of trajectories. As a result, slight shifts or scales in trajectories may considerably change the value of c_1 .

The mentioned interpretations can provide surgical mentors with good clues about the regular mistakes between beginners and help them to develop their lectures and training flowcharts accordingly. For instance, almost all novice and intermediate trainees can be confidently separated from expert ones according to their low value of c_1 in their embedding space shown in all plots of Fig. 6. In general, the low values of c_1 in the rotational data of the suturing task demonstrates the low angular dancing of the user while inserting the needle inside the tissue. This means the user mistakenly inserts the needle with translations rather than properly orienting the tool. This is mainly due to the bad needle positioning upon the surgical incision during the suture needle insertion. This is a common mistake in suturing tasks and this clue can be used as an important topic in training courses. To have a better understanding, please watch the supplementary video.

Moreover, bad needle positioning at the beginning of the task will propagate and lead to an extra Θ twist at the end. High values of c_3 in Fig. 6(e) demonstrates high pitch rotation at the end of suturing when the needle should come out of the tissue. It is the sign of inconvenient task completion due to the bad start that can harm the tissue and wrist of the surgeon.

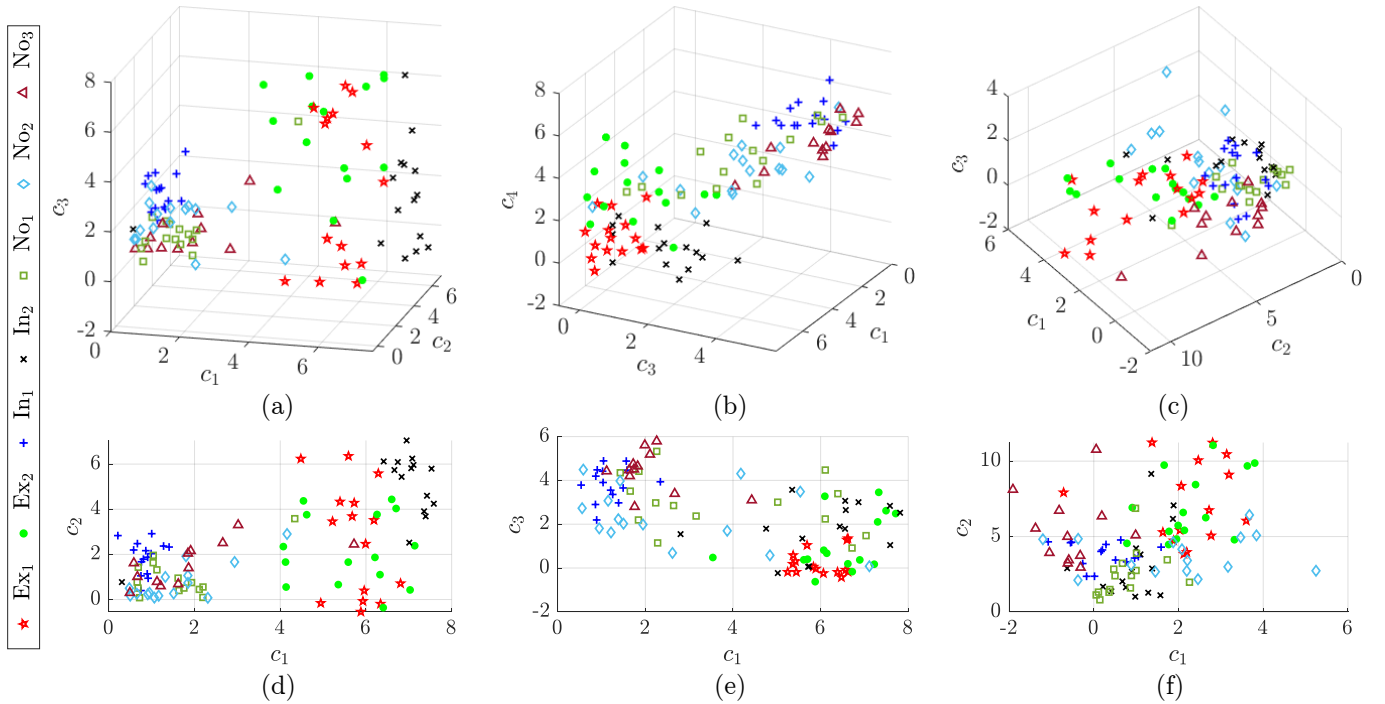


Figure 6: Embedding space representation of three different angular rotations (from the left column to right: Φ , Θ , and Ψ) of the suturing task for two expert users EX_i , two intermediate users IN_i , and three novice users NO_i trained over EX_1 trajectories in JIGSAWS data set. For better visualization, the second row shows 2D angle of the 3D plot of the same column.

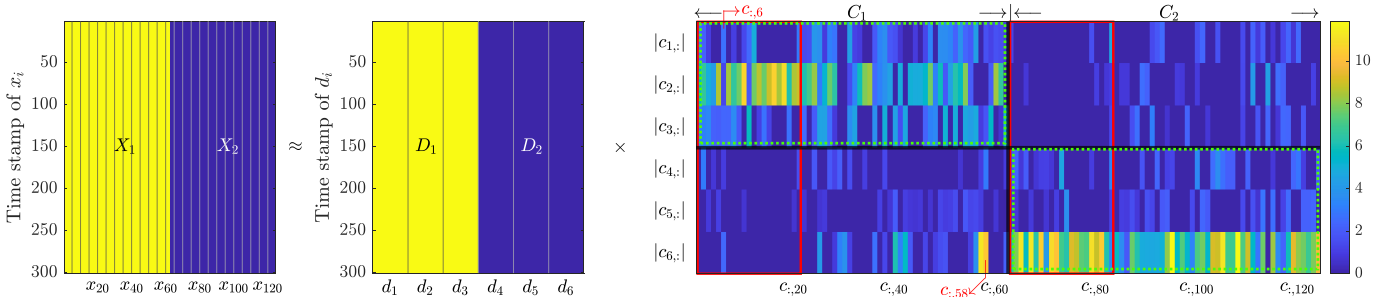


Figure 7: Bi-clustering of x and z trajectories of knot tying task in JIGSAWS data set according to the combined dictionary $D = [D_1, D_2]$ where D_i is trained over X_i . The generated code matrix $C = [C_1, C_2]$ for $X = [X_1, X_2]$ given D is roughly block-diagonal (dashed green rectangles). We can approximately assign each code vector into four main clusters according to the task (column clustering indicated by vertical black solid line) and user skills level (row clustering indicated by red rectangles). Red solid rectangles correspond to expert trials.

Additionally, since d_1^f in Fig. 2(d) is corresponding to the general trend in the trajectory, low values of c_1 in Fig. 6(e) can also demonstrate that less-experienced surgeons do not follow the general trend of the operation and the trajectory is mostly based on unordered motions.

E. Trajectory Bi-Clustering

Trajectory bi-clustering is another investigation that simultaneously reveals the *type* of the executive trajectory (e.g., x , y , z , Φ , Θ , and Ψ of suturing, knot tying, and needle passing tasks) and the skills level of the user. To do so, we concatenate different input matrices and their dictionaries to create $X = [X_1, \dots, X_t]$ and $D = [D_1, \dots, D_t]$ for t different input types where D_i is

trained over data set X_i . Then we let the sparse coding algorithm to generate code matrix C for X given the *combined* dictionary D . Ideally, different trajectories lie in non-overlapping subspaces and C generates codes for a given trajectory x_i based on its specifically designed trajectory and assigns zero to other non-related atoms. In other words, codes generated for a particular data type X_i is $C_i = [0_1^T, \dots, 0_{i-1}^T, D_i^T, 0_{i+1}^T, \dots, 0_t^T]^T$ where 0_l is a zero matrix with the same size as D_l for all $1 \leq l \leq t, l \neq i$. In this ideal case, the code matrix C is *block-diagonal*. Fig. 7 demonstrates the same concept for two types of inputs: x and z trajectories of knot tying task, in which the code matrix C is visualized by the colormap (i.e., warmer colors indicate higher code values.) The sparse coding algorithm tends to generate non-zero codes for a trajectory based on its specifi-

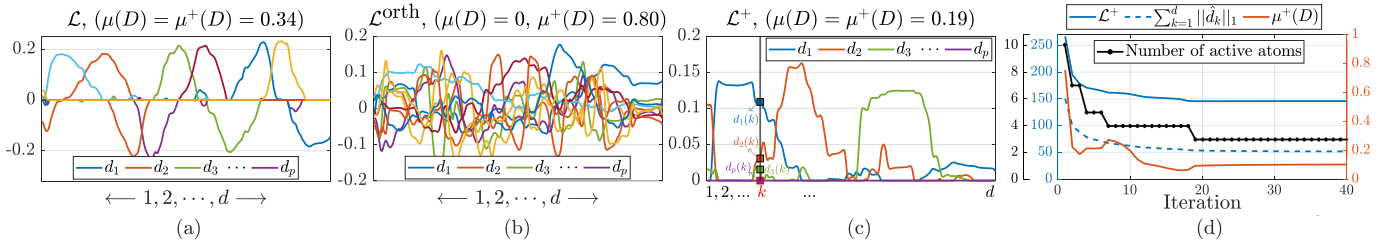


Figure 8: Dictionary atoms generated for a given task with $p=10$ initialized atoms based on (a) standard dictionary learning loss \mathcal{L} in (1), (b) orthogonal dictionary learning loss $\mathcal{L}^{\text{orth}}$ [24], and (c) our proposed \mathcal{L}^+ loss described in (4) (the algorithm nullified 7 atoms and preserved the most informative 3 atoms). (d) convergence properties of \mathcal{L}^+ loss (each y label is associated with the plot with the same color).

cally generated atoms which results in a block-diagonal code matrix \mathbf{C} (dashed green rectangles in Fig. 7). When the first dashed green block ends and the second one begins, we have a decision boundary between two clusters of input types (i.e., column clustering).

Moreover, expert users' trajectories have higher *fidelity* to their specific atoms. In other words, they purely belong to the subspace of their task with minimum overlap with the subspace of other tasks and almost do not generate codes for irrelevant atoms (i.e., their coefficients are zero). As a result, investigating the rows of the code matrix \mathbf{C} (i.e., row clustering) can give us informative clues about the skills level of the user (e.g., solid red rectangles in Fig. 7 which indicate trials done by expert users). In addition to the lack of expertise, abnormality is another cause of generating codes based on other irrelevant atoms which can be considered as random atoms for that specific trajectory. As it is shown in Fig. 7, 6th and 58th trials (i.e., $c_{:,6}$ and $c_{:,58}$) have small codes for their own atoms and large codes for other non-relevant atoms. It means for both trials, something wrong is happening during the execution of that trajectory that causes these trajectories to lose fidelity to their specifically generated atoms. Further explanations and the endoscopic videos of these two trials are provided in the supplementary [video](#).

V. DISCUSSIONS

A. Advantages of \mathcal{L}^+ Loss

In this section, we will compare final atoms generated by our method with \mathcal{L}^+ loss described in (4) with atoms generated by conventional methods with \mathcal{L} and $\mathcal{L}^{\text{orth}}$ losses to have a better understanding about the concepts behind these methods. As shown in Fig. 8(c), \mathcal{L}^+ generates optimal atoms compared to \mathcal{L} and $\mathcal{L}^{\text{orth}}$ in terms of minimum overlap, reduced $\mu(\mathbf{D})$ and $\mu^+(\mathbf{D})$ metrics, and good trajectory reconstruction. A result of the conventional dictionary learning problem is provided in Fig. 8(a). In general, conventional dictionary learning method with loss \mathcal{L} produces atoms with large overlaps and increased $\mu^+(\mathbf{D})$ compared to our approach since it does not have any penalization term for atoms' overlaps. Fig. 8(b) demonstrates atoms resulting from minimizing $\mathcal{L}^{\text{orth}}$ loss in which any pair of ten atoms have zero correlation (i.e., $\mu(\mathbf{D})=0$ or $\mathbf{D}^T\mathbf{D}=\mathbf{0}$), but they have considerable overlap with each other (i.e., large value for $\mu^+(\mathbf{D})$). It means orthogonal atoms generated by $\mathcal{L}^{\text{orth}}$ shown in Fig. 8(b), despite their low reconstruction loss and zero correlation cannot be used for approximate trajectory decomposition for structured tasks.

Table I: Ablation study on hyperparameters α and β in (4) for Φ_R .

α	β	Reconstruction loss	Final number of atoms	μ^+
0.8	1.4	9.07	6	0.19
0.85	1.4	9.86	5	0.20
0.9	1.4	10.59	4	0.21
1	1.5	11.37	3	0.10
1.7	1.5	13.56	2	0.23

Moreover, the final converged atoms based on \mathcal{L} and $\mathcal{L}^{\text{orth}}$ losses highly depend on the random state for the initialization. As illustrated in Section IV, our method is robust against initialization and delivers quite consistent results.

The effect of optimizing (4) on $\mu^+(\mathbf{D})$ is also illustrated in Fig. 8(d). The value of $\mu^+(\mathbf{D})$ drops when the sparsity promoting cost for dictionary atoms (i.e., $\sum_{k=1}^d \|\hat{d}_k\|_1$) decreases until the algorithm nullifies one or several atoms. When a reduction happens in the total number of active atoms, there is a mild and temporary increase in $\mu^+(\mathbf{D})$. This is because the algorithm is forced to assign data variations to fewer active atoms which increases the total overlap.

B. Hyperparameter Tuning

Another important feature of our method is that sparsity-promoting terms for both codes c_i and dictionary rows \hat{d}_k in \mathcal{L}^+ loss nullifies unimportant atoms and returns the fewer number of atoms after the optimization procedure. This fact is illustrated in Fig. 8(d) by showing the total number of active atoms in each iteration which is reducing when the optimization proceeds. As explained before, this feature benefits the ease of interpretability of generated atoms. By changing hyperparameters α and β , the number of final atoms and the value of their overlap will change. Ablation study results for right hand's Φ angle (Φ_R) in suturing task in Table I, indicate that relaxing the sparsity regularization parameters α and/or β yields an increased number of final atoms and an improvement in the reconstruction loss. However, arbitrarily increasing the number of atoms by reducing α and β does not guarantee a considerable reduction in reconstruction loss for unseen data and may lead the algorithm to become overfitted on the training set. As another perspective, trimming factors α and β iteratively reduce the degree-of-freedom of the algorithm and prevent the model from making redundant atoms. α and β can be empirically fine-tuned to reach the ideal number of minimally-overlapping atoms which is equal

to the intrinsic dimensionality of embedding space (δ). Although the normal range of hyperparameters α and β heavily depends on the application, for the examples of this work the empirical range is $0.1 \leq \alpha, \beta \leq 2.5$.

The hyperparameter γ regulates the relative importance of cross-correlation versus reconstruction loss in dictionary temporal fine-tuning post-training stage. Since capturing exact temporal position of each floating atom for each test sample has higher priority than reconstruction loss, empirically observed that $\gamma = 0.5$ is an acceptable choice for our experiments.

C. Final Number of Active Atoms

Determining the value of intrinsic dimensionality of embedding space (δ) needs field knowledge and depends on the task, the target variable for the investigation (e.g., translational data of the suturing task along x axis), and the number of independent sub-tasks within the given structured task (e.g., number of gestures in surgical trial). A good heuristic for finding this number is to plot the manifold of reconstruction error and μ^+ versus the number of active atoms and wherever the reconstruction error and/or μ^+ do not reduce with increasing the number of active atoms (i.e., the elbow of the manifold) we assign that number of atoms to δ . For instance, in Table I $\delta = 3$ is an ideal final number of atom since $\delta \geq 4$ suffers from large overlap between atoms (i.e., high value of μ^+) and $\delta = 2$ suffers from both poor reconstruction loss and large overlap between atoms. This result also makes intuitive sense; Φ_R trajectory in suturing task is composed of three main gestures: passing the needle inside the tissue, pulling the needle from the tissue, and passing the needle from one hand to the other one.

D. Rotational Data vs. Translational Data

Our observations suggest that, rotational data offer more interpretation and insight about the quality of executive tasks with sharper distinction between data clusters in the embedding space. One possible reason that rotational data are more *expressive* than translational data is that rotational patterns are more closely related to the skills level of the user. This is based on the intuition that humans according to their advanced motor-control capabilities, can accomplish a lot of complicated tasks by performing a succession of several motions, but the quality, dexterity, and efficiency will be determined based on how well they perform rotations while executing translations. As an example, bipedal robots exhibit the same problem. Although they follow human joint trajectories in walking task, the overall behavior of their walking is different from that of humankind.

Another important reason is the *curse of extra details* in translational data, which masks general patterns with unnecessary microscale motions and prevents meaningful atom generation. The reason might be due to minor mistakes users unconsciously correct with simple motions rather than sophisticated rotations. This makes the subject's behavior more streamlined in the rotation space and noisier in translation space. Finally, compared to other types of surgery, such as eye surgery, minimally invasive surgery offers many translations, which can increase the effects of random motions as well.

E. Applications in Bi-manual Tasks

In the proposed method, each trajectory component was investigated independently to analyze the performance of the executive task. However, in bimanual tasks (i.e., a class of tasks in which the brain must simultaneously plan and control the movements of both hands such as tying shoelaces or basic surgical tasks) what defines a person as an expert surgeon is not just simply what he/she performs by each individual hand but what he/she plans for the next step by executing a complex sequence of coordinated actions between his/her hands [33]. This concept is equivalent to the notion of *hands coordination* which measures the synchronicity and relationship between different trajectory components of one hand or between two hands. Incorporating coordination data together with the generated codes of the proposed method may benefit the quality of the final representation for downstream tasks such as skills classifier network. Extracting coordination data is beyond the scope of this paper and is a promising topic for future work.

VI. CONCLUSIONS

A new technique for visualizing surgical trajectories in the lowest possible dimensional space was presented in this paper. Representing trajectories based on a small number of minimally-overlapping atoms allowed us to meaningfully and intuitively decompose and investigate each trajectory. Each minimally-overlapping atom can be considered as a building block of the whole executive task that meaningfully reflects the style, skills, faults, and hidden behaviors of the user during performing the task. Incorporating floating atoms to capture important variations appearing at different temporal positions within the trajectory, improves model's accuracy and prevents information loss during the mapping procedure. All of these important features are objectivity expressed in terms of numerical gains in embedding (or code) space as an informative feature map for educational and examination purposes. According to our experiments on the JIGSAWS data set, our method is effective, reliable, and accurate for skills assessment and fault detection. In future work we plan to incorporate the idea of dynamic time warping (DTW) into our dual-sparse coding approach. We may find this helpful given the extra nonlinearity DTW introduces into our model.

APPENDIX

A. Proof of Property 1

According to *Hölder's inequality*, for any $\kappa, v \in (1, \infty)$ with $\frac{1}{\kappa} + \frac{1}{v} = 1$, we have

$$\sum_{k=1}^N |x_k y_k| \leq \left(\sum_{k=1}^N |x_k|^\kappa \right)^{\frac{1}{\kappa}} \left(\sum_{k=1}^N |y_k|^v \right)^{\frac{1}{v}} \quad (\text{A.1})$$

for all $\mathbf{x}^\top = (x_1, \dots, x_N)$, $\mathbf{y}^\top = (y_1, \dots, y_N) \in \mathbb{R}^N$. Considering the special case $\kappa = v = 2$ in Hölder's inequality, (A.1) for non-zero vectors \mathbf{x} and \mathbf{y} becomes

$$0 \leq |\mathbf{x}^\top \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \stackrel{\substack{\mathbf{x} \neq \mathbf{0} \\ \mathbf{y} \neq \mathbf{0}}}{=} 0 \leq \frac{|\mathbf{x}^\top \mathbf{y}|}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \leq 1. \quad (\text{A.2})$$

Applying (A.2) into (3) and considering the fact that $\mu(\mathbf{D}) \leq \mu^+(\mathbf{D})$ yields $0 \leq \mu(\mathbf{D}) \leq \mu^+(\mathbf{D}) \leq 1$ for any possible dictionary \mathbf{D} . ■

B. Convergence of the Coordinate Descent Algorithm 1

Consider the convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}).$$

Definition A.1 (η -strongly convex function): The following statements are all equivalent to the condition that a differentiable function \mathbf{f} is strongly convex with constant $\eta > 0$

- (i) $\mathbf{f}(\delta\mathbf{x} + (1-\delta)\mathbf{y}) \leq \delta\mathbf{f}(\mathbf{x}) + (1-\delta)\mathbf{f}(\mathbf{y}) - \frac{\delta(1-\delta)\eta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ for $\delta \in [0, 1]$.
- (ii) The function $\tilde{\mathbf{f}}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \frac{\eta}{2} \|\mathbf{x}\|_2^2$ is convex, $\forall \mathbf{x}$. \diamond

Definition A.2: The convex optimization objective function \mathbf{f} has component-wise L -Lipschitz continuous gradient if

$$\|\nabla_i \mathbf{f}(\mathbf{x} + h\mathbf{e}_i) - \nabla_i \mathbf{f}(\mathbf{x})\| \leq L|h| \quad (\text{A.3})$$

where $\mathbf{x} \in \mathbb{R}^N$, $h \in \mathbb{R}$, $i = 1, \dots, N$, and \mathbf{e}_i is the standard basis vector of i^{th} component in \mathbf{x} . \diamond

Lemma A.1: The criterion $\mathcal{L} = \sum_{i=1}^n \|\mathbf{x}_i - \sum_{j=1}^p \mathbf{d}_j c_{ij}\|_2^2$ equals the quadratic function $(\mathbf{D} - \mathbf{D}^*)^\top \mathbf{C}^\top \mathbf{C} (\mathbf{D} - \mathbf{D}^*)$ plus a constant where \mathbf{D}^* is the optimal solution of \mathbf{D} in minimizing \mathcal{L} . \square
Proof. See [28]. \blacksquare

Lemma A.2: The matrix $\mathbf{Q} = \mathbf{C}^\top \mathbf{C}$ defined in Lemma A.1 is positive definite. \square

Proof. At first, we will prove that \mathbf{Q} is a positive semidefinite matrix. According to the definition, the matrix \mathbf{Q} is positive semidefinite matrix if $\mathbf{z}^\top \mathbf{Q} \mathbf{z} \geq 0$, $\forall \mathbf{z} \in \mathbb{R}^N$. We have

$$\mathbf{z}^\top (\mathbf{C}^\top \mathbf{C}) \mathbf{z} = (\mathbf{C} \mathbf{z})^\top (\mathbf{C} \mathbf{z}) = \|\mathbf{C} \mathbf{z}\|_2^2 \geq 0.$$

Now, to prove that \mathbf{Q} is positive definite, we just need to prove that $\mathbf{z}^\top \mathbf{Q} \mathbf{z} \neq 0$, $\forall \mathbf{z} \neq \mathbf{0}$. Consider that $\mathbf{z}^\top \mathbf{Q} \mathbf{z} = 0$ for some $\mathbf{z} \neq \mathbf{0}$. It yields that $(\mathbf{D} - \mathbf{D}^*)^\top \mathbf{C}^\top \mathbf{C} (\mathbf{D} - \mathbf{D}^*)$ can be equal to zero for some $\mathbf{D} \neq \mathbf{D}^*$. It means, quadratic objective function \mathcal{L} has more than one optimal solution. Contradiction. As a result, $\mathbf{Q} = \mathbf{C}^\top \mathbf{C}$ is positive definite and $\lambda_{\min}(\mathbf{Q}) > 0$ where $\lambda_{\min}(\mathbf{Q})$ is the smallest eigenvalue of \mathbf{Q} . \blacksquare

Lemma A.3: The criterion \mathcal{L} defined in Lemma A.1 is η -strongly convex function with $\eta = 2\lambda_{\min}(\mathbf{C}^\top \mathbf{C}) > 0$. \square

Proof. According to Definition A.1(ii), $\tilde{\mathbf{f}}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \frac{\eta}{2} \|\mathbf{x}\|_2^2$ is convex if $\nabla^2 \tilde{\mathbf{f}}(\mathbf{x}) = \nabla^2 \mathbf{f}(\mathbf{x}) - \eta \mathbf{I} \succeq \mathbf{0}$ where $\mathbf{A} \succeq \mathbf{0}$ means that matrix \mathbf{A} is positive semidefinite. As a result, a twice continuously differentiable \mathbf{f} is η -strongly convex if $\nabla^2 \mathbf{f}(\mathbf{x}) \succeq \eta \mathbf{I}$, $\forall \mathbf{x}$ or equivalently, the smallest eigenvalue of $\nabla^2 \mathbf{f}(\mathbf{x})$ satisfies $\lambda_{\min}(\nabla^2 \mathbf{f}(\mathbf{x})) \geq \eta$, $\forall \mathbf{x}$. According to the quadratic form of the criterion \mathcal{L} , $\nabla^2 \mathbf{f}(\mathbf{x}) = \nabla^2 \mathcal{L} = 2\mathbf{Q}$. As a result, \mathcal{L} is η -strongly convex function with $\eta = 2\lambda_{\min}(\mathbf{C}^\top \mathbf{C})$ which according to Lemma A.2, $\eta > 0$. \blacksquare

Lemma A.4: Let $\mathcal{F} = \mathbf{f} + \mathbf{g}$ where \mathbf{f} is η -strongly convex function and \mathbf{g} is a convex function (not necessarily strongly convex), then \mathcal{F} is η -strongly convex function. \square

Proof. According to Definition A.1(i)

$$\mathcal{F}(\delta\mathbf{x} + (1-\delta)\mathbf{y}) = \mathbf{f}(\delta\mathbf{x} + (1-\delta)\mathbf{y}) + \mathbf{g}(\delta\mathbf{x} + (1-\delta)\mathbf{y})$$

$$* \leq \delta\mathbf{f}(\mathbf{x}) + (1-\delta)\mathbf{f}(\mathbf{y}) - \frac{\delta(1-\delta)\eta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \mathbf{g}(\delta\mathbf{x} + (1-\delta)\mathbf{y})$$

$$** \leq \delta\mathcal{F}(\mathbf{x}) + (1-\delta)\mathcal{F}(\mathbf{y}) - \frac{\delta(1-\delta)\eta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

where inequality * follows from the fact that \mathbf{f} is strongly convex and inequality ** holds since \mathbf{g} is convex. \blacksquare

Now, we want to prove the convergence of Algorithm 1 that tries to minimize (4). Both (5) and (10) are composed of a reconstruction error that according to Lemma A.1 equals to a strongly convex quadratic function plus a convex l_1 normalization part. As a result, according to Lemma A.4, (5) and (10) are strongly convex functions. In \mathbf{C} -step, we minimize (5) which according Lemma A.3 is strongly convex with $\eta_1 = 2\lambda_{\min}(\mathbf{C}^\top \mathbf{C}) > 0$. In \mathbf{D} -step we minimize (10) which according Lemma A.3 is strongly convex with $\eta_2 = 2\lambda_{\min}(\mathbf{D}\mathbf{D}^\top) > 0$. As a result, (4) is strongly convex with $\eta = \min\{\eta_1, \eta_2\} > 0$.

Lemma A.5: The cost function (5) is directional component-wise Lipschitz continuous gradient with $L_1 = \max_i \{\|\mathbf{d}_i\|_2^2\}$ in \mathbf{C} -step domain. \square

Proof. The gradient of (5) with respect to a parameter c_{ij} (i.e., the j^{th} component of a code vector \mathbf{c}_i) is

$$\nabla_{c_{ij}} \mathcal{L}_C^+ = \mathbf{d}_j^\top (\mathbf{D}\mathbf{c}_i - \mathbf{x}_i) + \alpha \text{sign}(c_{ij}). \quad (\text{A.4})$$

(A.4) is also used as a part of update rule in `sparseCode` algorithm to calculate the optimal value of code matrix \mathbf{C} . Due to the limitations rising from the learning rate and approximation nature of `sparseCode` algorithm, in cases that $c_{ij} \equiv 0$ the algorithm fails to land on $c_{ij} = 0$ and we have a fluctuation around the origin since $\text{sign}(c_{ij})$ fluctuates between -1 and 1 . A common heuristic applied in this situation is to *clamp* c_{ij} to zero (i.e., manually setting c_{ij} to zero) to prevent such fluctuations. Under these circumstances, we can assume that $\text{sign}(c_{ij})$ does not change and is the same for $\nabla_i \mathbf{f}(\mathbf{x} + h\mathbf{e}_i)$ and $\nabla_i \mathbf{f}(\mathbf{x})$ in Definition A.2. This assumption leads us to the directional component-wise L -Lipschitz continuous gradient in \mathbf{C} -step domain. According to (A.3) we should have

$$|\mathbf{d}_j^\top \mathbf{d}_j h \mathbf{e}_{ij}| = |\mathbf{d}_j^\top \mathbf{d}_j| |h| = \|\mathbf{d}_j\|_2^2 |h| \leq L|h|. \quad (\text{A.5})$$

We can reach to (A.5) for all components of code matrix \mathbf{C} . As a result, a good choice for L is the length of largest code vector \mathbf{d}_j in the dictionary matrix \mathbf{D} , or equivalently, $L_1 = \max_i \{\|\mathbf{d}_i\|_2^2\}$. \blacksquare

Corollary A.5.1: According to Lemma A.5, the cost function (10) is directional component-wise Lipschitz continuous gradient with $L_2 = \max_j \{\|\mathbf{c}_j\|_2^2\}$ in \mathbf{D} -step domain.

Corollary A.5.2: According to Lemma A.5 and Corollary A.5.1, the cost function (4) is directional component-wise Lipschitz continuous gradient with $L = \max\{L_1, L_2\}$ in the dual-sparse coding algorithm domain.

According to [29], [30] and the values of η and L calculated above, for each iteration $k \geq 0$ for Algorithm 1 starting from \mathbf{C}_0 and \mathbf{D}_0 , we have

$$\mathbb{E}[\mathcal{L}^+(\mathbf{D}_{k+1}, \mathbf{C}_{k+1})] - \mathcal{L}^+(\mathbf{D}^*, \mathbf{C}^*) \leq \left(1 - \frac{\eta}{2Ldp^2n}\right) [\mathcal{L}^+(\mathbf{D}_k, \mathbf{C}_k) - \mathcal{L}^+(\mathbf{D}^*, \mathbf{C}^*)] \quad (\text{A.6})$$

where \mathbf{D}^* and \mathbf{C}^* are local optima of (4) and d , p , and n are dimensions of matrices \mathbf{D} and \mathbf{C} defined in (1). Note that, if in Lemma A.5, our cost function was component-wise Lipschitz continuous gradient, not *directional* component-wise Lipschitz continuous gradient, \mathbf{D}^* and \mathbf{C}^* are global optimum solutions.

$2dp^2n$ is the total number of components in D and C which is the total number of parameters for the optimization problem (4). Although the values of η and L depend on the matrices D and C in the previous step, the convergence of expected error to zero (i.e., reaching to local or global optima) in each iteration is guaranteed in (A.6) since the value of $1 - \eta / (2Ldp^2n)$ is always less than 1.

REFERENCES

- [1] Abed Soleymani, Xingyu Li, and Mahdi Tavakoli. Deep neural skill assessment and transfer: Application to robotic surgery training. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021.
- [2] Abed Soleymani and et al. Surgical skill evaluation from robot-assisted surgery recordings. In *2021 International Symposium on Medical Robotics (ISMR)*. IEEE, 2021.
- [3] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Accurate and interpretable evaluation of surgical skills from kinematic data using fully convolutional neural networks. *International journal of computer assisted radiology and surgery*, 14(9):1611–1617, 2019.
- [4] Isabel Funke, Sören Torge Mees, Jürgen Weitz, and Stefanie Speidel. Video-based surgical skill assessment using 3d convolutional neural networks. *International journal of computer assisted radiology and surgery*, 14(7):1217–1225, 2019.
- [5] Hazel Doughty, Walterio Mayol-Cuevas, and Dima Damen. The pros and cons: Rank-aware temporal attention for skill determination in long videos. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7862–7871, 2019.
- [6] Niko Sünderhauf, Oliver Brock, Walter Scheirer, Raia Hadsell, Dieter Fox, Jürgen Leitner, Ben Upcroft, Pieter Abbeel, Wolfram Burgard, Michael Milford, et al. The limits and potentials of deep learning for robotics. *The International journal of robotics research*, 37(4-5):405–420, 2018.
- [7] Beatrice van Amsterdam, Matthew Clarkson, and Danaïl Stoyanov. Gesture recognition in robotic surgery: a review. *IEEE Transactions on Biomedical Engineering*, 2021.
- [8] Colin Lea, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks: A unified approach to action segmentation. In *European Conference on Computer Vision*, pages 47–54. Springer, 2016.
- [9] Giovanni Menegozzo, Diego Dall’Alba, Chiara Zandonà, and Paolo Fiorini. Surgical gesture recognition with time delay neural network based on kinematic data. In *2019 International Symposium on Medical Robotics (ISMR)*, pages 1–7. IEEE, 2019.
- [10] Robert DiPietro, Colin Lea, Anand Malpani, Narges Ahmidi, S Swaroop Vedula, Gyusung I Lee, Mija R Lee, and Gregory D Hager. Recognizing surgical activities with recurrent neural networks. In *International conference on medical image computing and computer-assisted intervention*, pages 551–558. Springer, 2016.
- [11] Robert DiPietro, Narges Ahmidi, Anand Malpani, Madeleine Waldram, Gyusung I Lee, Mija R Lee, S Swaroop Vedula, and Gregory D Hager. Segmenting and classifying activities in robot-assisted surgery with recurrent neural networks. *International journal of computer assisted radiology and surgery*, 14(11):2005–2020, 2019.
- [12] Danit Itzkovich, Yarden Sharon, Anthony Jarc, Yael Refaely, and Ilana Nisky. Using augmentation to improve the robustness to rotation of deep learning segmentation in robotic-assisted surgical data. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 5068–5075. IEEE, 2019.
- [13] Beatrice van Amsterdam, Matthew J Clarkson, and Danaïl Stoyanov. Multi-task recurrent neural network for surgical gesture recognition and progress prediction. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1380–1386. IEEE, 2020.
- [14] Daochang Liu and Tingting Jiang. Deep reinforcement learning for surgical gesture segmentation and classification. In *International conference on medical image computing and computer-assisted intervention*, pages 247–255. Springer, 2018.
- [15] Jacob Rosen, Massimiliano Solazzo, Blake Hannaford, and Mika Sinanan. Task decomposition of laparoscopic surgery for objective evaluation of surgical residents’ learning curve using hidden markov model. *Computer Aided Surgery*, 7(1):49–61, 2002.
- [16] Arthur Derathé, Fabian Reche, Pierre Jannin, Alexandre Moreau-Gaudry, Bernard Gibaud, and Sandrine Voros. Explaining a model predicting quality of surgical practice: a first presentation to and review by clinical experts. *International Journal of Computer Assisted Radiology and Surgery*, 16(11):2009–2019, 2021.
- [17] Ivana Tošić and Pascal Frossard. Dictionary learning. *IEEE Signal Processing Magazine*, 28(2):27–38, 2011.
- [18] Alexander L Chistov and D Yu Grigor’Ev. Complexity of quantifier elimination in the theory of algebraically closed fields. In *International Symposium on Mathematical Foundations of Computer Science*, pages 17–31. Springer, 1984.
- [19] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11(1), 2010.
- [20] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [21] Chenglong Bao, Hui Ji, Yuhui Quan, and Zuwei Shen. Dictionary learning for sparse coding: Algorithms and convergence analysis. *IEEE transactions on pattern analysis and machine intelligence*, 38(7):1356–1369, 2015.
- [22] Michael Elad and Michal Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing*, 15(12):3736–3745, 2006.
- [23] Joel A Tropp. Greed is good: Algorithmic results for sparse approximation. *IEEE Transactions on Information theory*, 50(10):2231–2242, 2004.
- [24] Chenglong Bao, Jian-Feng Cai, and Hui Ji. Fast sparsity-based orthogonal dictionary learning for image restoration. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3384–3391, 2013.
- [25] Michal Aharon, Michael Elad, and Alfred Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311–4322, 2006.
- [26] Carolyn Chen, et al. Crowd-sourced assessment of technical skills: a novel method to evaluate surgical performance. *Journal of surgical research*, 187(1):65–71, 2014.
- [27] Yixin Gao, S Swaroop Vedula, Carol E Reiley, Narges Ahmidi, Balakrishnan Varadarajan, Henry C Lin, Lingling Tao, Luca Zappella, Benjamin Béjar, David D Yuh, et al. Jhu-isi gesture and skill assessment working set (jigsaws): A surgical activity dataset for human motion modeling. In *Miccai workshop: M2cai*, volume 3, page 3, 2014.
- [28] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [29] Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [30] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1):1–38, 2014.
- [31] Timothy N Judkins, Dmitry Oleynikov, and Nick Stergiou. Objective evaluation of expert and novice performance during robotic surgical training tasks. *Surgical endoscopy*, 23(3):590–597, 2009.
- [32] Ignacio Ramirez, Pablo Sprechmann, and Guillermo Sapiro. Classification and clustering via dictionary learning with structured incoherence and shared features. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3501–3508. IEEE, 2010.
- [33] Vedula, S Swaroop and Malpani, Anand and Ahmidi, Narges and Khudanpur, Sanjeev and Hager, Gregory and Chen, Chi Chung Grace. Task-level vs. segment-level quantitative metrics for surgical skill assessment. *Journal of surgical education*, 73(3):482–489, 2016.