

# Deep learning-enabled prediction of surgical errors during cataract surgery: from simulation to real-world application

Maxime Faure<sup>a,b</sup>, Pierre-Henri Conze<sup>a,c</sup>, Béatrice Cochener<sup>a,b,d</sup>, Anas-Alexis Benyoussef<sup>a,b,d</sup>, Mathieu Lamard<sup>a,b</sup>, Gwenolé Quellec<sup>a</sup>

<sup>a</sup>LaTIM UMR 1101, Inserm, Brest, France

<sup>b</sup>University of Western Brittany, Brest, France

<sup>c</sup>IMT Atlantique, Brest, France

<sup>d</sup>Ophthalmology Department, University Hospital of Brest, Brest, France

## Abstract

Real-time prediction of technical errors from cataract surgical videos can be highly beneficial, particularly for telementoring, which involves remote guidance and mentoring through digital platforms. However, the rarity of surgical errors makes their detection and analysis challenging using artificial intelligence. To tackle this issue, we leveraged videos from the EyeSi Surgical cataract surgery simulator to learn to predict errors and transfer the acquired knowledge to real-world surgical contexts. By employing deep learning models, we demonstrated the feasibility of making real-time predictions using simulator data with a very short temporal history, enabling on-the-fly computations. We then transferred these insights to real-world settings through unsupervised domain adaptation, without relying on labeled videos from real surgeries for training, which are limited. This was achieved by aligning video clips from the simulator with real-world footage and pre-training the models using pretext tasks on both simulated and real surgical data. For a 1-second prediction window on the simulator, we achieved an overall AUC of 0.820 for error prediction using 600×600 pixel images, and 0.784 using smaller 299×299 pixel images. In real-world settings, we obtained an AUC of up to 0.663 with domain adaptation, marking an improvement over direct model application without adaptation, which yielded an AUC of 0.578. To our knowledge, this is the first work to address the tasks of learning surgical error prediction on a simulator using video data only and transferring this knowledge to real-world cataract surgery.

**Keywords:** cataract surgery, capsulorhexis, surgical error prediction, deep learning, real-time video analysis, unsupervised domain adaptation

## 1. Introduction

### 1.1. Context

Cataract, defined as the opacification of the crystalline lens, is a prevalent ocular condition that significantly impacts vision worldwide [1]. Cataract surgery is the most performed surgical procedure worldwide, highlighting its critical role in ophthalmology and public health [2]. The large volume of procedures has generated extensive data, offering opportunities to improve quality management, education, and training [3]. In practice, phacoemulsification is the most effective method for lens removal [4], involving the replacement of the clouded lens with an artificial one. Successful execution of the capsulorhexis step, a crucial stage in cataract surgery, is critical to prevent complications. This technique involves creating a circular opening in the anterior lens capsule to access and remove the cloudy lens. For optimal results, the capsulorhexis must be circular, regular, and well-centered, without radial extensions or capsular tags which can lead to surgical complications [5]. Advances in artificial intelligence (AI), particularly deep learning, have demonstrated potential in evaluating surgeon performance through

video analysis [6]. These technologies have been shown to objectively assess surgical skills and support intraoperative tasks such as real-time tool identification [7, 8], and surgical phase recognition [9, 10, 11]. Video analysis offers benefits beyond the operating room, such as remote surgical supervision and telementoring, and has shown promise in detecting specific errors, including those related to intraocular lens (IOL) implantation [12], continuous curvilinear capsulorhexis and phacoemulsification [13]. However, there remains a lack of comprehensive data on surgical errors. While public datasets like CaDIS [14], 101-Cataract [15], and Cataract-1K [16] exist, the scarcity of annotated data for surgical error detection remains a major challenge. Surgical simulators have emerged as valuable tools to address this limitation and generate substantial data for analysis. These simulators provide a controlled environment for collecting extensive datasets, enabling the study of surgical techniques and error detection. Various simulation methods are employed in cataract surgery training, including virtual reality, wet-lab, dry-lab models, and e-learning for technical and non-technical skills [17]. Among these, the EyeSi simulator is the most widely used [18]. The effectiveness of the EyeSi simulator in ophthalmology training has been widely demonstrated. It accurately replicates many aspects of cataract surgery, provid-

Email address: maxime.faure@univ-brest.fr (Maxime Faure)

ing essential training for residents [19]. Simulator-based training has also been shown to reduce operative times for surgical residents learning phacoemulsification, compared to traditional methods [20]. A strong link between simulator performance and real-life outcomes has been established. [21] found that the number of surgeries performed by surgeons — reflecting their practical experience — correlated with simulator scores. Studies have shown correlations between simulator scores and various metrics, such as GRASIS [22] scores [23], motion-tracking performance [24], and OSACSS [25] scores, effectively distinguishing between novice and experienced surgeons [26]. Moreover, the use of simulators has significantly reduced postoperative complication rates. Inexperienced surgeons had a 27.14% complication rate, compared to 12.86% for those with intermediate experience, with the difference being statistically significant [27]. Capsular rupture rates also decreased by 38% for surgical residents using the simulator [28], and capsulorhexis errors were reduced by 68% after the simulator was introduced to training programs [29].

### 1.2. Prediction of Surgical Errors

A significant portion of research on surgical error detection focuses on gesture error detection, often utilizing robotic data, particularly kinematic data [30], and occasionally relying on image data [31]. Beyond robotics, most image-based studies prioritize real-time detection of adverse events in surgical videos, with bleeding being a primary target [32]. In the specific case of detecting surgical errors in cataract surgery, one study focuses on capsulorhexis-related events [13]. The authors demonstrated that the system, trained to detect surgical errors in real-time, was capable of making predictions in advance. However, these predictions did not accurately localize the errors in the future, which made it impossible to generate relevant alerts. Despite these advancements, there is a noticeable gap in research on predicting future surgical errors, which is essential for generating timely alerts and recommendations. To address this gap, this work aims to position itself within the field of event anticipation in videos, a domain known as action anticipation [33].

Action anticipation in deep learning typically involves analyzing temporal video segments to predict future actions or events. Video-based models can be used, as well as approaches that incorporate spatial and temporal encoders. While some studies forecast future events to support classification tasks using temporally conditioned models [34], others enhance predictions by incorporating additional data modalities, such as optical flow [35] or domain-specific knowledge about relevant objects [36]. There are numerous applications, although few are found in the field of surgery, and they primarily focus on surgical phase anticipation [37].

### 1.3. Domain Transfer and Adaptation

An effective approach for domain adaptation is data transformation, such as histogram matching [38], which facilitates data homogenization. Furthermore, adversarial techniques, including Domain-Adversarial Neural Networks (DANN) [39]

and their specialized adaptations for video data [40], as well as attention-based methods [41], demonstrate considerable potential. However, these adversarial methods often necessitate large batch sizes to perform optimally, as noted by [42]. Statistical discrepancy-based techniques, such as CORrelation ALignment (CORAL) [43] and Maximum Mean Discrepancy (MMD) [44], focus on aligning feature distributions by minimizing statistical differences between domains. These methods are particularly effective for feature alignment and are adaptable for video data. We have selected these approaches as the foundation for the development of our proposed algorithm.

#### 1.3.1. Objectives of the study

This work aims to demonstrate the feasibility of predicting surgical errors during capsulorhexis in real life.

We first developed two novel datasets to enable the prediction of surgical errors during capsulorhexis: the first consists of annotated capsulorhexis videos from a surgical simulator, capturing a range of error types, while the second is a dataset of real surgery videos, where a subset was annotated. The unannotated videos were used for unsupervised domain adaptation. These original datasets serve as the foundation for training and evaluating our models. Building on this, we implemented a real-time deep learning system that bridges the simulator’s final evaluation with real-time surgical error prediction. The simulator provides a variety of surgical error scenarios often absent in real videos, making it especially valuable for training a deep learning algorithm. Additionally, we demonstrated that knowledge could be transferred from the simulator to real-world scenarios through unsupervised domain adaptation. Our approach leveraged the extensive data from the simulator and real surgical videos without requiring labor-intensive annotation to train an algorithm that generates real-time alerts.

To our knowledge, this is the first attempt to transfer knowledge from simulators to real-world cataract surgeries.

## 2. Dataset and annotation

This section presents the acquisition of videos from the simulator and real surgeries, their preparation, preprocessing, and annotation.

### 2.1. Data

#### 2.1.1. Simulator data

The capsulorhexis module of the EyeSi® Surgical cataract surgery simulator allows the operator to practice the capsulorhexis step of the cataract surgery.

A total of 422 exercises were collected. These were performed by 11 users of the simulator located at Brest University Hospital between May 2022 and May 2023.

The average duration for this exercise is 106.8 seconds, the standard deviation is 56.9 seconds, ranging from 25 to 222 seconds.

For each exercise, we recorded video using a single camera view (monocular), capturing frames every 33ms (at a 30Hz rate). Their spatial resolution is 800×600 pixels.

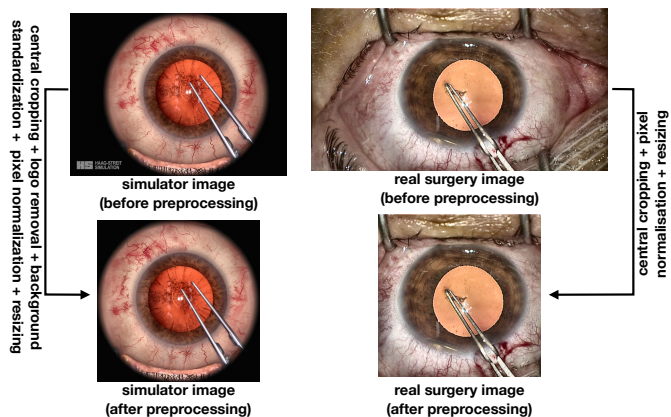


Figure 1: Steps and illustration of preprocessing results for images from simulator videos (left) and real surgery (right).

For video preprocessing, we cropped each image around the eye and ensured that all pixels representing the ocular border mask had a value of 0 through pixel normalization. The image size was then adjusted to 600×600 pixels. The preprocessing result is shown in Fig. 1.

Note that we also have access to the eye’s three Euler angles across time,  $\theta_x$ ,  $\theta_y$ , and  $\theta_z$ , which are used for annotating some surgical errors. These angles are expressed in a fixed reference frame and sampled at 30Hz.

### 2.1.2. Real surgery data

Between June and September 2023, we collected 107 monocular cataract surgery videos performed by five ophthalmologists at Brest University Hospital. Two surgeons had performed over 1000 surgeries, while the other three had completed fewer than 1000. We have a single video for each patient.

We excluded 21 of these videos as they were either unusable (incomplete capsulorhexis, eye not in the camera frame) or depicted scenarios not represented in the simulator videos (e.g., white cataract).

The capsulorhexis step, our focus here, was manually extracted from those videos. In the following sections, when we refer to the available real surgical videos, we mean only the videos showing the capsulorhexis step.

The average duration for performing capsulorhexis is 36.3 seconds, with a standard deviation of 20.4 seconds and ranging from 10 to 97 seconds. The videos are sampled at 30Hz, with a spatial resolution of 1280×720. For preprocessing, the video images were normalized so that pixel values fell within the range of 0 to 1 and the image size was adjusted to 600×600 pixels (Fig. 1).

## 2.2. Surgical errors

### 2.2.1. Description of simulator surgical errors at the video level

The surgical simulator assesses performance and provides a score sheet at the end of the exercise, highlighting surgical errors in the video. The surgical errors evaluated by the simulator include:

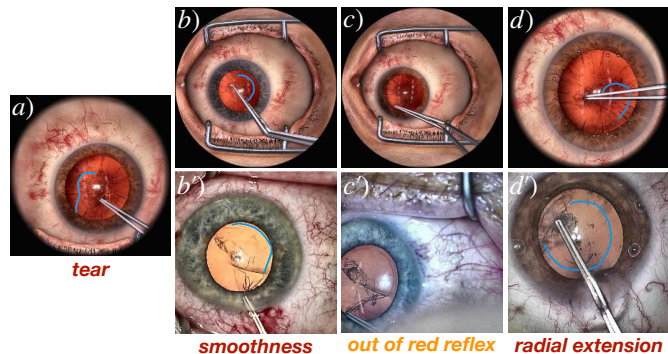


Figure 2: Illustrations of surgical errors on simulator video images and real surgical image: tear (a), smoothness (b, b’), out of red reflex (c, c’), and radial extension (d, d’). The capsulorhexis opening completed up to this point is highlighted in blue.

	simulator		real surgery	
	number of surgical error	duration (s)	number of surgical error	duration (s)
<i>smoothness</i>	760	×	27	×
<i>out of red reflex</i>	756	712	46	45
<i>radial extension</i>	226	453	9	6
<i>tear</i>	44	×	0	×

Table 1: Number of surgical errors and duration (in seconds) for the 422 annotated simulator videos and the 29 annotated real surgical videos.

- *operating without red reflex*: time spent with the eye out of red reflex. This metric is solely dependent on the eye and is not specific to the capsulorhexis step of cataract surgery.
- *smoothness*: roundness of the capsulorhexis, expressed as a percentage,
- *radial extension*: maximum radius of the capsulorhexis in millimeters.
- *tear*: the presence of a capsular tear, meaning a capsule rupture.

It is important to note that the occurrence of a surgical error does not necessarily indicate a failure; a capsulorhexis can still be deemed successful even if it is not perfect. Some surgical errors are less critical than others. For example, 100% of capsulorhexis cases with a tear are considered failures, while this is not the case for other errors.

### 2.2.2. Annotation of surgical errors over time

An expert proceeded with the localization of surgical errors within the videos during a data annotation phase at the frame level to know exactly when surgical errors occur.

For annotation, we set aside videos that do not contain surgical errors according to the score sheets and annotated the remaining ones only.

Here are the error classes we defined, illustrated in Fig. 2 (a, b, c, d, and e):

- The *out of red reflex* (Fig. 2c) error is automatically annotated using the eye’s Euler angles at time  $t$  if the empirically verified condition  $\theta_{xt}^2 + \theta_{yt}^2 > \left(\frac{\pi}{12}\right)^2$  is met.

- The initial frames at the start of a loss of circularity/regularity are annotated as containing a *smoothness* error (Fig.2b)
- When the capsulorhexis flap is too far from the center (Fig.2d): *radial extension*.
- If a tear is present on an image (Fig.2a): *tear*

Regarding the statistics from the simulator, among the 422 videos, *red reflex* errors occurred 756 times with a median duration of 0.67 seconds. *Radial extension* errors occurred 226 times (1.73 seconds), *tear* errors only 44 times, and *smoothness* errors occurred 760 times.

Fig.3 illustrates the surgical errors co-occurrences matrix  $M$  between  $err_A$  and  $err_B$  in simulator data:

$$M(err_A, err_B) = 100 \times \frac{|S_A \cap S_B|}{|S_A|}. \quad (1)$$

where  $S_A$  and  $S_B$  represent sequences of binary labels of length  $T$  (total duration of simulator videos in second), with  $S_A, S_B \in \{0, 1\}^T \times \{0, 1\}^T$  indicate the presence (1) or absence (0) of an error at each 1-second time sequence.  $M(err_A, err_B)$  denote the proportion of the error sequence  $S_A$  that coincides with the errors in  $S_B$ .

The co-occurrence of errors is generally low, as the occurrence of one error does not necessarily lead to another, except for *tear*  $\rightarrow$  *radial extension*. The asymmetry of the matrix shows that the reverse is not true.

While we do not have score sheets for the real surgical videos, annotations are performed on a subset of 29 videos. Surgical errors in real surgeries are much less frequent — approximately 23 times fewer.

The least observed surgical error in the simulator videos (*tear*) is absent in annotated real surgeries, as shown in Tab.1. The median duration for the *red reflex* error is 0.73 seconds, while it is 0.70 seconds for the *radial extension*. Some annotations are illustrated in Fig.2 (b', c', and d').

### 2.3. Data split

Simulator videos are split into 60/20/20% for the training/validation/test sets. We ensured that the surgical error content was evenly distributed across each set.

The split for the dataset of real surgical videos is 60/28/12%. This was done on unannotated videos, ensuring each split contained the same proportion of videos recorded on the same day. The split was necessary to account for slight visual and skill variations in the videos due to different camera settings and the surgeon performing the procedure. This ensured that the diversity of video appearances and surgeon expertise was evenly distributed across the datasets.

The 29 validation and test videos of the real video dataset were annotated afterward.

The training, validation, and test sets for simulator and surgical videos are decomposed into video snippets for error prediction.

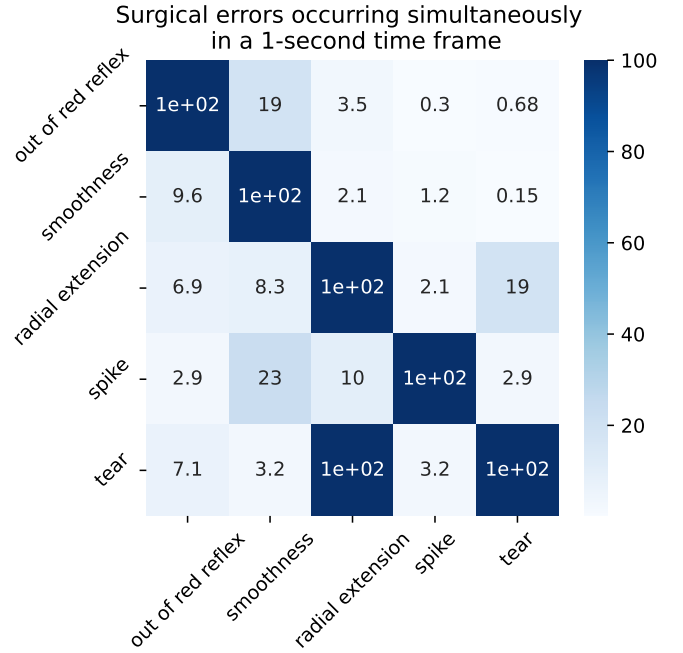


Figure 3: Co-occurrence of surgical errors within a 1-second window on simulator data.

## 3. Methods

### 3.1. Decomposition into video snippets

Simulator and surgical videos are segmented into snippets for the training, validation, and test sets. To simplify, we did not include the names of the videos from which each snippet is derived in the notation. Instead, each snippet is indexed by  $k$ .

Let  $A_{[t-D_k:t]}^{(k)}$  denote the prior video snippet at time  $t$ , lasting  $D_k$  seconds, and composed of  $L_k$  frames. For labeled videos, we also define  $P_{[t+\Delta T_k:t+\Delta T_k+1]}^{(k)}$  as the posterior sequence, starting  $\Delta T_k$  seconds (the prediction horizon) after the end of  $A_{[t-D_k:t]}^{(k)}$  and lasting for 1 second. We denote  $S_k \in \mathbb{N}^*$  such that  $S_k - 1$  represents the number of frames skipped between two consecutive frames in the original video. This parameter determines the temporal downsampling applied to create the video snippet. Since the base sampling on the simulator is 30 frames per second, the temporal duration  $D_k$  of  $A_{[t-D_k:t]}^{(k)}$  in seconds is:

$$D_k = \frac{L_k \times S_k}{30} \quad (2)$$

Thus, by choosing  $L_k$  and  $S_k$ , we control the extent of the history to predict surgical errors. For example, by choosing  $L_k = 30$  and  $S_k = 1$  or  $L_k = 10$  and  $S_k = 3$ , we have  $D_k = 1$  s. The prior video snippet is observed for prediction by the algorithm, while the posterior sequence is used to define the label (surgical error class)  $I_{\Delta T_k}$ , where  $I_{\Delta T_k} \in \{0, 1\}^{C+1}$ . In our study,  $C = 5$  represents the number of surgical classes, with the addition of +1 corresponding to one additional class for all surgical errors combined. Specifically, we aim to predict from

	number of snippet	proportion (%)
<i>smoothness</i>	18,787	21
<i>out of red reflex</i>	10,361	12
<i>radial extension</i>	5,443	6.2
<i>tear</i>	294	0.34
<i>no error</i>	50,893	58.8
total	86,468	100

Table 2: Distribution of surgical errors at the video snippet scale for the training dataset generated with  $L_k = 10$ ,  $S_k = 3$  and  $\Delta T_k = 1$  (simulator data).

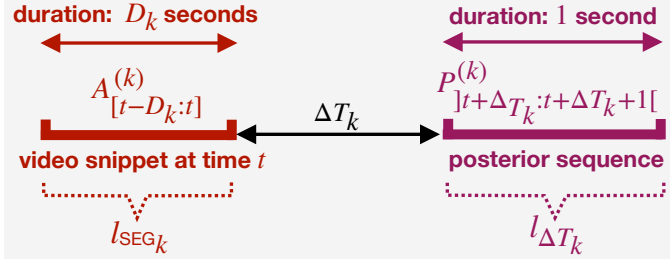


Figure 4: Illustration of the concepts composing our dataset:  $A_{[t-D_k:t]}^{(k)}$  the video snippet observed by the model of duration  $D_k$  seconds, which defines the current label  $I_{SEG_k}$ , and the posterior sequence  $P_{[t+\Delta T_k:t+\Delta T_k+1]}^{(k)}$ , of duration 1 second used to define the prediction label  $I_{\Delta T_k}$ . Note that the data is indexed by  $k$  and that the original video notation is not involved.

the video snippet  $A_{[t-D_k:t]}^{(k)}$  whether a surgical error will occur across the sequence  $P_{[t+\Delta T_k:t+\Delta T_k+1]}^{(k)}$ .

Concretely:

$$I_{\Delta T_k}[j] = \begin{cases} 1 & \text{if a surgical error of class } j \text{ occurs} \\ & \text{across the sequence } P_{[t+\Delta T_k:t+\Delta T_k+1]}^{(k)} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

We also classify the current error in the observed video snippet by defining the current label  $I_{SEG_k} \in \{0, 1\}^{L_k, C+1}$ . This label corresponds to the surgical error content in the observed video snippet, ensuring that the extracted features during training focus more on the current error information rather than other irrelevant details (Fig.4). We set  $\Delta T_k = T$  with  $T \in [1, 5]$  seconds. We arbitrarily shifted the sequences by 1/15 second for simulator data to create a dataset with diverse sequences. Additionally, we used markers for the training set and validation set to select every third clip without surgical errors, which helped balance the dataset. For the surgical videos, which are fewer in number, we consistently shifted by 1/30. Tab.2 illustrates the number of video snippets associated with each type of surgical error for a training dataset constructed with  $L_k = 10$  and  $S_k = 3$ . Note that we ensured to have the same number of video snippets and the same proportion of surgical errors for the considered configurations  $(L_k, S_k)$  to ensure fair comparisons between various settings.

### 3.2. Surgical error prediction from simulator data

In this subsection, we present the prediction model that will be trained in a supervised manner using data from the simulator.

#### 3.2.1. Algorithm for surgical error prediction

We aim to implement a real-time analysis system. Therefore, we develop a conditional algorithm to make predictions with a customizable prediction horizon.

We denote by  $\mathcal{F}$  the spatial encoder model,  $\mathcal{M}$  the temporal encoder model,  $C^{(c)}$  two classification layers for the current ( $c$ ) classification,  $C^{(p)}$  two classification layers for prediction ( $p$ ), att an attention layer and TE two fully-connected layers.

Thus, we have:

- $SR_k = \mathcal{F}(A_{[t-D_k:t]}^{(k)})$ . It corresponds to applying  $\mathcal{F}$  to each element of the sequence  $A_{[t-D_k:t]}^{(k)}$ .
- $TR_k = \mathcal{M}(SR_k)$
- $CC_k = C^{(c)}(TR_k)$  and  $CC_k^{(att)} = \hat{I}_{SEG_k} = \text{att}(CC_k)$
- $TR_k^{(+)} = \text{concat}(TR_k, CC_k^{(att)}, \text{TE}(\Delta T_k))$
- $\hat{I}_{\Delta T_k} = FP_k = C^{(p)}(TR_k^{(+)})$

with  $SR_k \in \mathbb{R}^{1, L_k, I}$  being the spatial representation with  $I$  as the latent representation size of the images.

$TR_k \in \mathbb{R}^{1, S}$  is the temporal representation of the data with  $S$  being the chosen size of this latent representation.  $CC_k \in \mathbb{R}^{1, L_k, C+1}$  (current classification) is the surgical error classification for each image of the observed video snippet.  $CC_k^{(att)} \in \mathbb{R}^{1, C+1}$  then corresponds to the average surgical error (averaged by an attention mechanism) of the previous sequence. Thus,  $TR_k$  and  $CC_k$  can be seen as information relative to the previous sequence.

$\text{TE}(\Delta T_k) \in \mathbb{R}^{1, T}$  corresponds to a learned temporal representation of the prediction horizon  $\Delta T_k$ .

$TR_k^{(+)} \in \mathbb{R}^{1, S+C+1+T}$  is the concatenation of the temporal representation, the predicted surgical error content of the observed sequence, and the prediction horizon to customize the future prediction  $FP_k \in \mathbb{R}^{1, C+1}$ . We trained a conditional model to avoid training multiple models for predictions across different prediction horizons  $\Delta T_k$ .

The upper part of Fig.5 shows the different parts of the deep learning algorithm used for error prediction in surgical procedures based on the simulator.

Note that for the sake of clarity, we have fixed the batch size to  $B = 1$  here.

#### 3.2.2. Loss Function

To train the model, we minimize the empirical risk defined by a cross-entropy loss function with two terms:

$$\begin{aligned} \mathcal{L}_C(.) &= \lambda_1 \mathcal{L}_C^{(c)}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(c)}}) \\ &\quad + \mathcal{L}_C^{(p)}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(c)}}, W_{\text{TE}}, W_{C^{(p)}}) \\ &= -\frac{\lambda_1}{B} \frac{1}{B} \sum_{k=1}^B \sum_{c=1}^C I_{SEG_k}[c] \log(\hat{I}_{SEG_k}[c]) - \sum_{k=1}^B \sum_{c=1}^C I_{\Delta T_k}[c] \log(\hat{I}_{\Delta T_k}[c]) \end{aligned} \quad (4)$$

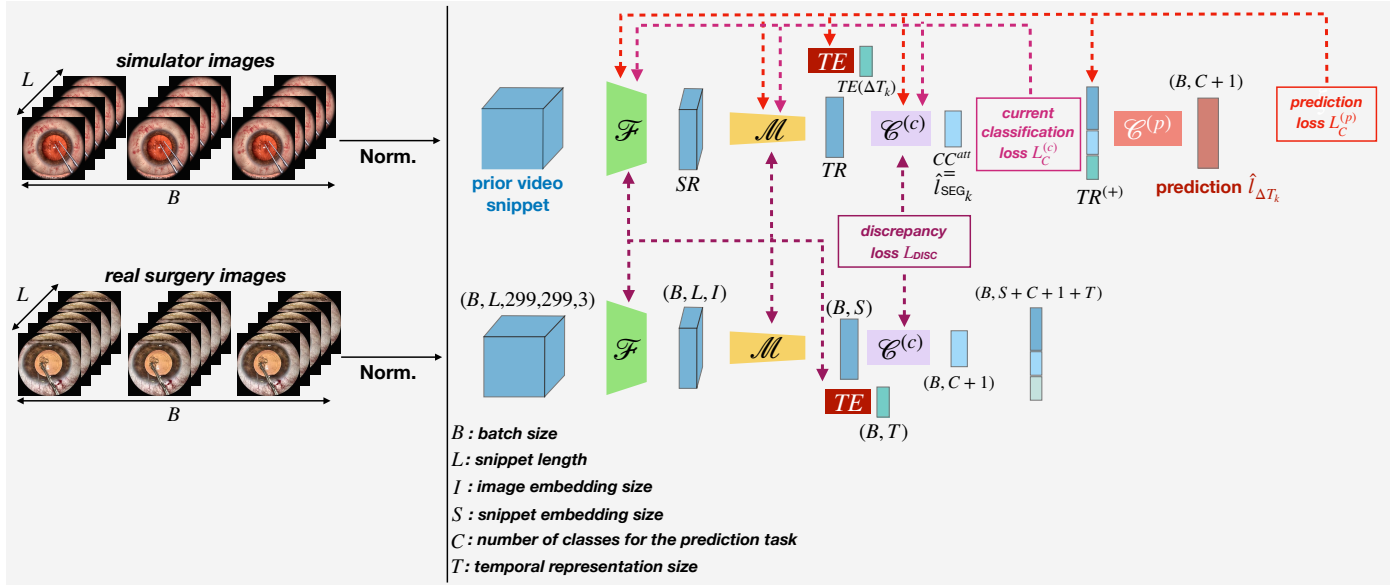


Figure 5: Simplified diagram of the unsupervised domain adaptation approach based on distribution alignment. The dotted arrows represent the backpropagation of the gradient.

The first term  $\mathcal{L}_C^{(c)}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{\mathcal{C}^{(c)}})$  is the current classification loss, where  $W_{\mathcal{F}}$  corresponds to the learned weights of the model  $\mathcal{F}$ . It is associated with the classification of the surgical error that occurred in the previous sequence. This term enables the model to learn the observed sequence’s surgical error content and generate a richer representation to guide the final classification. The second term corresponds to the prediction made within the prediction horizon  $\Delta T_k$ .

### 3.3. Unsupervised domain adaptation

After addressing the prediction of surgical errors using simulator videos, the focus now shifts to adapting to real surgical videos without available annotations for training. This challenge necessitates unsupervised domain adaptation techniques to predict surgical errors in real videos by leveraging knowledge from annotated simulator data while adapting the differences between the two domains.

#### 3.3.1. Data homogenization

Homogenizing the data is the most straightforward strategy when considering knowledge transfer from one domain to another. The goal is to directly bring the target domain data closer to the source domain.

We applied several preprocessing steps to the images from real surgeries, as illustrated in Fig.6.

- *central crop*: Extracting a central square crop from the image. This is the standard preprocessing technique mentioned in Sect.2.1.1.
- *ocular border mask*: Adding black pixels around the central square crop to align the histograms of the images (simulator images contain many black pixels).

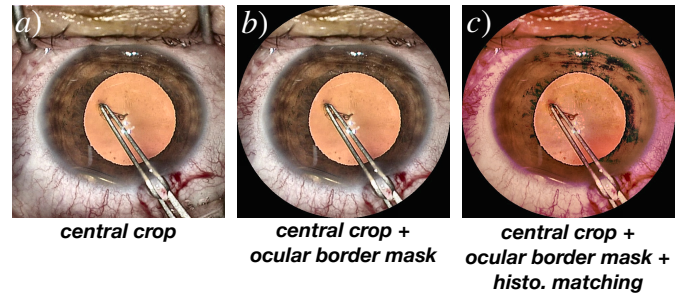


Figure 6: Illustrations of images from real surgical videos: central crop from a simulator image (a), images with ocular border mask (b), and images with ocular border mask and histogram matching (c).

- *ocular border mask + histogram matching*: Histogram matching is an image processing technique used to adjust the pixel intensity distribution of one image to match the histogram of another [38]. This method is beneficial to standardize brightness, contrast, or color distribution between images. The intuition behind applying this technique is to align the color tones of real surgical video frames with those of the simulator, reducing variance and simplifying the task for the network. Several reference images (10) were selected based on iris color variability and magnification level to ensure a straightforward and effective computation.

This strategy does not include any additional training. A model trained on simulator data will be directly evaluated on the transformed video snippets from the real surgery validation set.

### 3.3.2. Feature alignment

As the variance of the images at the pixel level is very large, we have considered a transfer approach that passes through the latent space.

This is an unsupervised approach that does not exploit the labels of the classification on the simulator.

The idea is to align the representations of the video snippets from the simulator using a discrepancy loss. We chose to perform this alignment with a statistical method that brings the distributions closer by considering a batch of video snippets.

In our study, we have considered the CORrelation ALignment (CORAL) method [43], which aims to align the covariance matrices of the extracted features for the source and target data by minimizing their difference.

$$\mathcal{L}_{\text{CORAL}}(\cdot) = \frac{1}{4S^2} \|C_S - C_T\|_F^2 \quad (5)$$

where  $\|\dots\|_F^2$  represents the Frobenius norm and with:

$$C_S = \frac{1}{B-1} \left( \mathbf{TR}_{(S)}^T \mathbf{TR}_{(S)} - \frac{1}{B} (\mathbf{1}^T \mathbf{TR}_{(S)})^T (\mathbf{1}^T \mathbf{TR}_{(S)}) \right) \quad (6)$$

$$C_T = \frac{1}{B-1} \left( \mathbf{TR}_{(T)}^T \mathbf{TR}_{(T)} - \frac{1}{B} (\mathbf{1}^T \mathbf{TR}_{(T)})^T (\mathbf{1}^T \mathbf{TR}_{(T)}) \right) \quad (7)$$

the covariance matrices estimated from the batch of  $B$  features of the video snippets from the source (simulator) ( $S$ ) and target (real surgery) ( $T$ ) domains, respectively.

This loss can be seen as an alignment constraint. In parallel, the classification task is performed solely on the simulator data.

We also considered in our study the Maximum Mean Discrepancy (MMD), which is a measure between latent representations involving a kernel [44].

$$\begin{aligned} \mathcal{L}_{\text{MMD}}^2(\cdot) &= \frac{1}{B(B-1)} \sum_i \sum_{j \neq i} \varphi(\mathbf{TR}_{(S)}[i], \mathbf{TR}_{(S)}[j]) \\ &\quad - 2 \frac{1}{B^2} \sum_i \sum_j \varphi(\mathbf{TR}_{(S)}[i], \mathbf{TR}_{(T)}[j]) \\ &\quad + \frac{1}{B(B-1)} \sum_i \sum_{j \neq i} \varphi(\mathbf{TR}_{(T)}[i], \mathbf{TR}_{(T)}[j]) \end{aligned}$$

$$\text{with } \varphi(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{-\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}\right).$$

We also compared these losses to the Mean Squared Error (MSE):

$$\mathcal{L}_{\text{MSE}}(\cdot) = \frac{1}{B} \sum_i (\mathbf{TR}_{(S)}[i] - \mathbf{TR}_{(T)}[i])^2$$

The overall loss is therefore:

$$\mathcal{L}(\cdot) = \mathcal{L}_C(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(p)}}, W_{\text{TE}}) + \lambda_2 \mathcal{L}_{\text{DISC}}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(p)}}, W_{\text{TE}}, W_{C^{(p)}})$$

where we set  $\mathcal{L}_{\text{DISC}}$  as either  $\mathcal{L}_{\text{CORAL}}$ ,  $\mathcal{L}_{\text{MMD}}$ , or  $\mathcal{L}_{\text{MSE}}$ . Fig.5 illustrates the training procedure for feature alignment.

### 3.3.3. Pre-training

Self-supervised learning is particularly useful in domain adaptation when dealing with multiple domains, as it enables the definition of tasks that extract feature sets common to both domains. This approach can facilitate subsequent classification tasks and potentially enhance performance [45].

We defined two self-supervised learning tasks, also known as pretext tasks, to pre-train the model.

These tasks are performed simultaneously on both the simulator data and real surgery data, allowing the model to learn shared representations from both domains and improve its ability to adapt to domain-specific challenges.

- **Frame Order Task:** This task is formulated as a binary classification problem where the model predicts whether the frames in the input video snippet are in the correct temporal order. The associated cross-entropy loss is  $\mathcal{L}_C^{(p)}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(p)}})$  with  $C^{(p)}$  represents the classification layer for this pretext task.
- **Image Reconstruction Task:** This task is framed as an image reconstruction problem. The model is trained to reconstruct a randomly cropped quarter of an image after a prediction horizon  $\Delta T_k$ . The objective is to minimize the mean squared error (MSE) between the predicted and target images. The loss function is defined as  $\mathcal{L}_{\text{MSE}}(W_{\mathcal{F}}, W_{\mathcal{M}}, W_{\text{att}}, W_{C^{(p)}}, W_{\text{TE}}, W_{\mathcal{D}})$  guiding the model to accurately predict the image across the temporal gap. The decoder,  $\mathcal{D}$ , is an LSTM-CNN symmetric to the encoder.

## 4. Experiments and results

### 4.1. Evaluation metric

The Area Under the ROC Curve (AUC) is chosen as the evaluation metric for this study.

### 4.2. Training settings

We implemented the algorithms using the PyTorch framework and utilized the Adam optimizer with a learning rate of  $5 \times 10^{-5}$  for the spatial encoder  $5 \times 10^{-4}$  for the rest.

For training the model to predict surgical errors using only the simulator data, we employed batches of size  $B = 32$  snippets to minimize empirical risk with gradient accumulation. We set  $\lambda_1 = 0.1$  (see 3.2.2) and trained for 10 epochs, employing early stopping based on overall AUC to prevent overfitting.

For domain adaptation, and computational efficiency, we used a batch size of  $B = 8$  and trained for 10 epochs with  $\lambda_2 = 1$  (see 3.3.2).

Given the significant variability in the results from validation sets, we conducted five separate training runs. We selected the best outcomes based on their performance on the validation set.

### 4.3. Surgical error prediction on simulator data

For this study, we simplified the setup by setting the previous sequence to a duration of  $D_k = 1$  second, with  $L_k = 10$  frames, corresponding to a sampling rate of  $S_k = 3$  samples/s.

Unless otherwise specified, the prediction horizon was set to 1 second for training and inference.

We evaluated multiple spatial and temporal encoders. Tab.3 presents the comparative study results. The inference time shown is the average time in milliseconds to compute  $SR_k$  over 1000 inferences. An asterisk (\*) denotes spatial encoders pre-trained on images of different sizes than those used in our study (fourth column).

For temporal encoders, we used the 3 best configurations without pre-training:

- LSTM: 2 layers, 256 hidden units,
- 1D CNN: dilated causal convolution [46],
- Transformer encoder: 3 layers, 8 attention heads, and a 256-dimensional embedding.

In our comparative study, 2D CNN models surpassed vision Transformers (with spatial encoder) and 3D CNN model [47]. The best-performing spatial encoders were:

- Inception\_ResNetv2 pre-trained on ImageNet using 299×299 pixel images, further trained for surgical error prediction using 600×600 pixel images
- EfficientNetB0 pre-trained on noisy student using 224×224 pixel images, further trained for surgical error prediction using 600×600 pixel images

However, EfficientNetB0 achieved the highest AUC of 0.812 with 600×600 pixel images on the validation data, with statistical significance (p-value = 1.89e-05, DeLong test [48]).

In general, the LSTM outperformed other temporal encoders or achieved similar results, except for the Transformer when using EfficientNetB7 (AUC = 0.795 compared to 0.784, with a p-value of 1.55e-05), and InceptionV3 with an image size of 299×299 (AUC = 0.773 compared to 0.766, p-value of 1.77e-05).

In all cases, the results are better when error prediction is performed using 600×600 images compared to smaller sizes (224×224, 299×299, 384×384 pixels), even when the model was pre-trained with smaller image sizes. This underscores the importance of larger image sizes for capturing comprehensive eye movement and information related to the capsulorhexis flap.

Note that the best result with a size smaller than 600×600 pixels was achieved using Inception\_ResNetv2 and LSTM, which was pre-trained on ImageNet with a size of 299×299 pixels and trained for surgical error prediction with 299×299 pixels, resulting in an AUC of 0.784. Tab.4 presents the results from the second comparative study. We maintained the best configuration (EfficientNetB0 + LSTM, 600×600 pixel images) and varied sequence hyperparameters ( $L_k$  and  $S_k$ ). The prediction horizon remained fixed at 1 second for inference. Note

that, here, the inference time corresponds to the average time to compute future predictions for a batch size of 1.

The study shows that using a single input image yield unsatisfactory results (AUC = 0.734).

Better results were achieved with a 1/3-second window or a 1-second window using  $(L_k, S_k) = (10, 3)$  or  $(L_k, S_k) = (5, 6)$ , yielding significantly better results than video snippets with 30 frames  $(L_k, S_k) = (30, 1)$  and reduced computation time by around 25 %.

Tab.5 provides results by surgical error classes with EfficientNetB0 + LSTM, 600×600 pixel,  $(L_k, S_k) = (10, 3)$  images for the same prediction horizon ( $\Delta T_k = 1$  second) on validation and test data.

Fig.7 shows the performance across different prediction horizons  $\Delta T_k \in [1, 5]$  for 3 different training strategies.

The first strategy involves training one model for each prediction horizon. The second strategy consists of training a single model sequentially, starting with  $\Delta T_k = 1$  and incrementally progressing up to  $\Delta T_k = 5$ . The third strategy entails this sequential training without providing the temporal horizon as input to the model during the training process.

### 4.4. Surgical error prediction transfer to real surgical videos

Due to limitations in computation time and memory capacity, we consider Inception\_ResNetv2 and use images of size 299×299 pixels for unsupervised domain adaptation. We set  $L_k = 10$ ,  $S_k : 3$  and  $\Delta T_k = 1$  second.

Tab.6 presents the unsupervised domain adaptation results, using data homogenization and feature alignment strategies on validation dataset. The prediction horizon is fixed at 1 second for inference. We provide the overall results and those for the three specific surgical error classes observed in real surgical videos.

The *central crop* of the real surgical images in video snippets to match the simulator’s format served as the baseline, representing the lower performance bound for transfer.

Adding an *ocular border mask* improved performance across all surgical error classes (+6.5% overall).

The *histogram matching* strategy does not seem to have provided any benefit, except for the *radial extension* class (+18%), but the AUC of 0.509 indicates a performance only slightly better than random.

Regarding feature alignment, the row labeled *source* corresponds to the validation results for the simulator data, while the row labeled *real surgery target* indicates the results on the real surgery validation set. Note that the symbol  $\emptyset$  represents the results obtained without using the discrepancy loss.

Tab.7 presents results for feature alignment using the CORAL discrepancy loss after pre-training the model with 2 different pretext tasks on validation and test data.

## 5. Discussion

This study aimed to predict real-time surgical errors from monocular cataract surgical videos by leveraging knowledge transferred from simulator data, with only a limited amount of



Spatial encoder	Version	Number of parameters (M)	Image size	FLOPs (G)	Average inference time (ms)	Temporal encoder		
						LSTM	1D-CNN	Transformer
ResNet	50	25.5	224	4	$7.2 \pm 0.04$	$0.749 \pm 0.006$	$0.746 \pm 0.005$	$0.749 \pm 0.006$
	50*		600	29	$7.5 \pm 0.06$	$0.789 \pm 0.005$	$0.768 \pm 0.007$	$0.759 \pm 0.007$
Inception_ResNet	v2	56.4	299	13	$30.1 \pm 2.2$	$0.784 \pm 0.005$	$0.754 \pm 0.007$	$0.755 \pm 0.006$
	v2*		600	56	$42.4 \pm 0.03$	$0.796 \pm 0.005$	<b><math>0.783 \pm 0.007</math></b>	<b><math>0.802 \pm 0.005</math></b>
EfficientNet	b0ns	5	224	0.4	$8.1 \pm 1.9$	$0.776 \pm 0.006$	$0.743 \pm 0.005$	$0.732 \pm 0.006$
	b0ns*		600	3	$9.6 \pm 2.5$	<b><math>0.812 \pm 0.005</math></b>	$0.768 \pm 0.006$	$0.780 \pm 0.005$
	b7ns	66	600	38	$54.6 \pm 0.05$	$0.784 \pm 0.006$	$0.774 \pm 0.005$	$0.795 \pm 0.005$
InceptionNet	v3	24	299	6	$12.2 \pm 1.2$	$0.766 \pm 0.006$	$0.751 \pm 0.005$	$0.773 \pm 0.006$
	v3*		600	24	$13.6 \pm 1.7$	$0.789 \pm 0.005$	$0.751 \pm 0.007$	$0.753 \pm 0.006$
ViT	vit_base_patch16_224	87	224	18	$10.8 \pm 1.8$	$0.752 \pm 0.006$	$0.731 \pm 0.007$	$0.734 \pm 0.006$
	vit_base_patch16_384	87	384	56	$12.8 \pm 1.7$	$0.779 \pm 0.006$	$0.758 \pm 0.005$	$0.764 \pm 0.006$
						None		
X3D	M	4	224	5	$19 \pm 1.3$	$0.762 \pm 0.008$		

Table 3: Overall results (AUC) on the validation simulator dataset for various temporal and spatial encoders, with  $L_k = 10$ ,  $S_k = 3$ , and  $\Delta T_k = 1$  s. Results are presented as the Mean (95% CI), derived from the DeLong test. Significant results are highlighted in bold.

$D_k$ (s)	$L_k$	$S_k$	inference time (ms)	AUC (overall)
1/30	1	1	$6.8 \pm 0.6$	$0.734 \pm 0.006$
1/3	5	2	$8.9 \pm 1.9$	<b><math>0.815 \pm 0.005</math></b>
	10	1	$9.6 \pm 0.9$	<b><math>0.820 \pm 0.006</math></b>
1	5	6	$8.9 \pm 2.2$	<b><math>0.811 \pm 0.005</math></b>
	10	3	$9.9 \pm 1.9$	<b><math>0.812 \pm 0.005</math></b>
	30	1	$13.4 \pm 2.5$	$0.785 \pm 0.006$
2	10	6	$9.6 \pm 1.6$	$0.799 \pm 0.005$
	30	2	$12.0 \pm 2.3$	$0.791 \pm 0.006$
	60	1	$16.8 \pm 2.7$	$0.782 \pm 0.006$

Table 4: Overall results (AUC) on the validation simulator dataset obtained with the best models for various sequence hyperparameters ( $L_k$ ,  $S_k$ ) and  $\Delta T_k = 1$  s. Significant results are highlighted in bold.

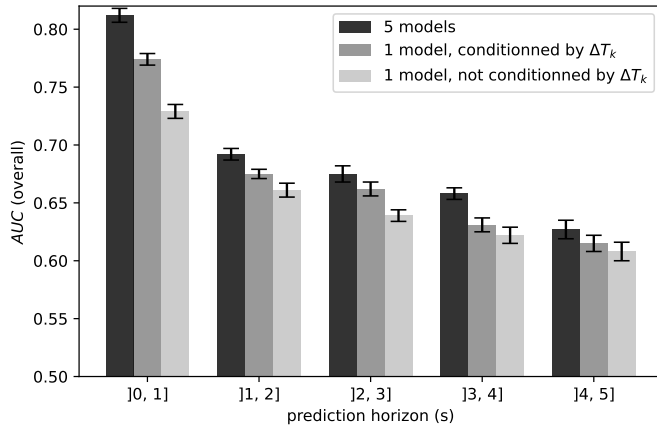


Figure 7: Validation results (AUC) on the simulator dataset, obtained with the best configuration (algorithm and hyperparameters) for various prediction horizons.

labeled data from actual surgical videos. The use of a simulator offers a wide range of scenarios, allowing for more comprehensive predictions.

As expected, prediction becomes more challenging as the prediction horizon increases (e.g.,  $AUC = 0.812$  for  $\Delta T_k = 1$  second versus  $0.692$  for  $\Delta T_k = 2$  seconds). The use of a time-step model yields significantly better results for prediction horizons  $\Delta T_k$  ranging from 1 to 5 seconds. However, it has been shown that if a single model is employed, conditional predictions result in reduced performance decline for the shorter prediction horizons ( $\Delta T_k = 1, 2,$  and  $3$  seconds). Our results suggest that having an extensive temporal history is unnecessary for accurate surgical error prediction. A short history of just 1 second proved to be sufficient, which is advantageous for real-time analysis from a computational point of view. Moreover, downsampling using every second frame did not degrade performance, enabling us to process video snippets with fewer frames and thus accelerating computations, which is beneficial for real-time analysis. However, an analysis based solely on individual images seems insufficient, incorporating at least a basic level of temporal dynamics is essential. According to our comparative study, using  $600 \times 600$  pixel images resulted in significantly better overall results. The best performance was achieved with EfficientNetB0, pre-trained on noisy student as spatial encoder, and an LSTM as temporal encoder. This configuration yielded an AUC score of  $0.820$  for a 1-second prediction horizon on the test set with  $L_k = 10$  and  $S_k = 3$ . Regarding surgical error classes, we found that not all surgical errors are equally predictable. For instance, the *out of red reflex* error is easier to predict as it depends on the eye's orientation, i.e. general geometric features the model can readily capture.

Regarding the transfer to real surgery videos, the results on real surgical videos are promising, validating this initial attempt of knowledge transfer for the prediction task. We observe that the homogenization strategy, specifically adding an ocular border mask, improves performance compared to the baseline

	AUC (overall)	AUC (smooth.)	AUC (out of red reflex)	AUC (radial extension)	AUC (tear)
<i>validation</i>	0.812 ± 0.005	0.789 ± 0.006	0.976 ± 0.003	0.698 ± 0.013	0.804 ± 0.026
<i>test</i>	0.820 ± 0.005	0.771 ± 0.006	0.981 ± 0.003	0.682 ± 0.011	0.815 ± 0.025

Table 5: Prediction results (AUC) for surgical errors by class, obtained with the best model on validation and test simulator data.

		AUC (overall)	AUC (smooth.)	AUC (out of red reflex)	AUC (radial extension)	
<i>data homogenization</i>	<i>central crop</i>	0.570	0.485	0.649	0.430	
	<i>ocular border mask</i>	0.615	0.559	0.678	0.480	
	<i>ocular border mask + histo. matching</i>	0.561	0.502	0.639	0.509	
<i>features alignment</i>	<i>simulator (source)</i>	$\emptyset$	<b>0.784</b>	<b>0.715</b>	0.950	<b>0.647</b>
		<i>CORAL</i>	0.771	0.645	<b>0.963</b>	0.593
		<i>MMD</i>	0.747	0.700	<b>0.967</b>	0.602
		<i>MSE</i>	0.765	<b>0.719</b>	0.956	0.617
	<i>real surgery (target)</i>	<i>CORAL</i>	0.654	0.611	0.710	0.751
		<i>MMD</i>	0.651	0.605	0.685	0.629
		<i>MSE</i>	0.542	0.517	0.645	0.518

Table 6: AUC results (domain adaptation) for data homogenization and feature alignment. Significant results are highlighted in bold.

	AUC (overall)	AUC (smooth.)	AUC (out of red reflex)	AUC (radial extension)
<i>Frame order</i>	0.663	0.632	0.757	0.658
<i>Image reconstruction</i>	0.657	0.617	0.774	0.674

Table 7: AUC results for feature alignment using CORAL discrepancy loss with pre-training on the validation real surgery dataset.

(AUC = 0.615 vs. 0.578 overall). This mask enhances the structural consistency of the images by adding a black pixel border. Simple histogram matching, however, does not appear to be effective, yielding lower performance than the ocular border mask alone (AUC = 0.561 vs. 0.578 overall). This may be due to the overly strict constraint of histogram matching and the global, rather than local, changes in light intensity. The differences between the domains are not only explained by luminescence but also by the structure’s shape, which is not altered by functions that act on the histogram. We also observe that aligning latent representations using CORAL and MMD yields better performance than data homogenization alone (AUC = 0.654 and 0.651 overall vs 0.615). The results obtained with these discrepancy losses are better than those with MME, which does not allow unsupervised domain adaptation in this work. It is important to note the slight drop in performance on the source domain compared to results without latent representation alignment (AUC = 0.771 with CORAL, 0.747 with MMD vs 0.784 overall). This decline is expected, as aligning latent representations involves a trade-off between optimizing surgical error prediction and regularizing the latent space. However, the degradation is minimal, highlighting that the transfer does not come at the cost of accurate surgical error prediction. Surprisingly, better performance (+1.5%) was achieved for the out-of-red reflex error when applying domain adaptation compared to when no adaptation ( $\emptyset$ ) was performed. This is likely due to the training strategy—using five models for transfer learning and retaining the best results, compared to using one without transfer—rather than the effects of the discrepancy loss. The temporal pretext task followed by domain adaptation with the CORAL loss improves overall performance (+2%) and for *smoothness* (around

+3%), but especially for the *out of red reflex* error, with an improvement of nearly 7%. On the other hand, the image reconstruction-based pretext task shows an advantage mainly for the *out of red reflex* error, with a 9% increase. This is not surprising, as this error class is easier to transfer because it relies primarily on geometric features, which are more straightforward to generalize across domains.

While our approach focused on data alignment, future research could explore alternative strategies such as data augmentation. Data augmentation could help addressing both inter-domain and intra-domain variability, potentially leading to further performance improvements. Finally, future research could explore integrating advanced information beyond RGB channels, such as detailed information about the capsulorhexis process itself (identifying the capsulorhexis opening). Additionally, the use of binocular video could further enhance surgical error prediction. Binocular videos provide richer data by capturing depth and spatial relationships more effectively, which can facilitate more accurate analysis of surgical movements and errors. This additional perspective could give a more comprehensive view of the surgical field, potentially enhancing model performance and the reliability of predictions.

A major limitation is the very low number of errors in real surgeries, which makes validation and testing challenging, and comparison to a supervised domain adaptation approach impossible. Moreover, although the simulator contains many errors compared to a real surgery dataset, the classes are highly imbalanced. Novel strategies should be established to address this challenge.

## Compliance with ethical standards

This study was conducted following the Declaration of Helsinki. Consent was obtained from participants to allow this study to take place.

## Acknowledgments

The authors thank Haag-Streit for making this work possible by sharing the simulator data and authorizing it. We also wish to express our gratitude to the ophthalmology department of the University Hospital of Brest for their collaboration and for providing essential data for our study. This work was supported by Inserm and the Brittany Region through the ARED program (ASCIA project).

## References

- [1] T. Rossi, M. R. Romano, D. Iannetta, V. Romano, L. Gualdi, I. D'Agostino, G. Ripandelli, Cataract surgery practice patterns worldwide: a survey, *BMJ Open Ophthalmology* 6 (2021) e000464. doi:10.1136/bmjophth-2020-000464.
- [2] Eurostat, Surgical operations and procedures statistics, [https://ec.europa.eu/eurostat/statistics-explained/index.php?title=Surgical\\_operations\\_and\\_procedures\\_statistics](https://ec.europa.eu/eurostat/statistics-explained/index.php?title=Surgical_operations_and_procedures_statistics), 2023. Accessed: [Insert access date here].
- [3] S. Müller, M. Jain, B. Sachdeva, P. Shah, F. Holz, R. Finger, K. Murali, M. Wintergerst, T. Schultz, Artificial intelligence in cataract surgery: A systematic review, *Translational vision science & technology* 13 (2024) 20. doi:10.1167/tvst.13.4.20.
- [4] G. Davis, The evolution of cataract surgery, *Missouri Medicine* 113 (2016) 58–62.
- [5] H. V. Gimbel, T. Neuhann, Continuous curvilinear capsulorhexis, *Journal of Cataract and Refractive Surgery* 17 (1991) 110–111. doi:10.1016/S0886-3350(13)81001-2.
- [6] R. Rampat, R. Deshmukh, X. Chen, D. S. W. Ting, D. G. Said, H. S. Dua, D. S. J. Ting, Artificial intelligence in cornea, refractive surgery, and cataract: Basic principles, clinical applications, and future directions, *Asia-Pacific Journal of Ophthalmology* 10 (2021) 268–281. doi:10.1097/APO.0000000000000394.
- [7] H. Al Hajj, M. Lamard, P.-H. Conze, B. Cochener, G. Quellec, Monitoring tool usage in surgery videos using boosted convolutional and recurrent neural networks, *Medical Image Analysis* 47 (2018) 203–218. URL: <https://www.sciencedirect.com/science/article/pii/S1361841518302470>. doi:https://doi.org/10.1016/j.media.2018.05.001.
- [8] N. Matton, A. Qalieh, Y. Zhang, A. Annadanam, A. Thibodeau, T. Li, A. Shankar, S. Armenti, S. I. Mian, B. Tannen, N. Nallasamy, Analysis of cataract surgery instrument identification performance of convolutional and recurrent neural network ensembles leveraging bigcat, *Translational Vision Science & Technology* 11 (2022) 1. doi:10.1167/tvst.11.4.1.
- [9] G. Quellec, K. Charrière, M. Lamard, Z. Droueche, C. Roux, B. Cochener, G. Cazuguel, Real-time recognition of surgical tasks in eye surgery videos, *Medical Image Analysis* 18 (2014) 579–590. URL: <https://www.sciencedirect.com/science/article/pii/S1361841514000309>. doi:https://doi.org/10.1016/j.media.2014.02.007.
- [10] G. Quellec, M. Lamard, B. Cochener, G. Cazuguel, Real-time task recognition in cataract surgery videos using adaptive spatiotemporal polynomials, *IEEE Transactions on Medical Imaging* 34 (2015) 877–887. doi:10.1109/TMI.2014.2366726. arXiv:2014 Oct 31.
- [11] S. Touma, F. Antaki, R. Duval, Development of a code-free machine learning model for the classification of cataract surgery phases, *Scientific Reports* 12 (2022) 2398. URL: <https://doi.org/10.1038/s41598-022-06127-5>. doi:10.1038/s41598-022-06127-5.
- [12] N. Ghamsarian, M. Taschwer, D. Putzgruber-Adamitsch, S. Sarny, Y. El-Shabrawi, K. Schoeffmann, LensID: a CNN-RNN-based framework towards lens irregularity detection in cataract surgery videos, in: *Medical Image Computing and Computer Assisted Intervention*, 2021, pp. 76–86.
- [13] S. Morita, H. Tabuchi, H. Masumoto, H. Tanabe, N. Kamiura, Real-time surgical problem detection and instrument tracking in cataract surgery, *Journal of Clinical Medicine* 9 (2020) 3896. doi:10.3390/jcm9123896.
- [14] M. Grammatikopoulou, E. Flouty, A. Kadkhodamohammadi, G. Quellec, A. Chow, J. Nehme, I. Luengo, D. Stoyanov, Cadis: Cataract dataset for surgical rgb-image segmentation, *Medical Image Analysis* 71 (2021) 102053. URL: <https://www.sciencedirect.com/science/article/pii/S1361841521000992>. doi:https://doi.org/10.1016/j.media.2021.102053.
- [15] K. Schoeffmann, M. Taschwer, S. Sarny, B. Münzer, M. J. Primus, D. Putzgruber, Cataract-101: video dataset of 101 cataract surgeries, in: P. César, M. Zink, N. Murray (Eds.), *Proceedings of the 9th ACM Multimedia Systems Conference, MMSys 2018, Amsterdam, The Netherlands, June 12–15, 2018, ACM*, 2018, pp. 421–425. URL: <https://doi.org/10.1145/3204949.3208137>. doi:10.1145/3204949.3208137.
- [16] N. Ghamsarian, Y. El-Shabrawi, S. Nasirihaghghi, D. Putzgruber-Adamitsch, M. Zinkernagel, S. Wolf, K. Schoeffmann, R. Sznitman, Cataract-1k dataset for deep-learning-assisted analysis of cataract surgery videos, *Scientific Data* 11 (2024) 373. URL: <https://doi.org/10.1038/s41597-024-03193-4>. doi:10.1038/s41597-024-03193-4.
- [17] R. Lee, N. Raison, W. Y. Lau, A. Aydin, P. Dasgupta, K. Ahmed, S. Halidar, A systematic review of simulation-based training tools for technical and non-technical skills in ophthalmology, *Eye* 34 (2020) 1737–1759. URL: <https://doi.org/10.1038/s41433-020-0832-1>. doi:10.1038/s41433-020-0832-1.
- [18] C. Cissé, K. Angioi, A. Luc, J.-P. Berrod, J.-B. Conart, Eyesi surgical simulator: validity evidence of the vitreoretinal modules, *Acta Ophthalmologica* 97 (2019) e277–e282. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/aos.13910>. doi:https://doi.org/10.1111/aos.13910. arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1111/aos.13910.
- [19] L. Carr, T. McKechnie, A. Hatamnejad, J. Chan, A. Beattie, Effectiveness of the eyesi surgical simulator for ophthalmology trainees: systematic review and meta-analysis, *Canadian Journal of Ophthalmology* 59 (2024) 172–180. URL: <https://www.sciencedirect.com/science/article/pii/S0008418223001072>. doi:https://doi.org/10.1016/j.jcjo.2023.03.014.
- [20] T. M. Ahmed, B. Hussain, M. A. R. Siddiqui, Can simulators be applied to improve cataract surgery training: a systematic review, *BMJ Open Ophthalmol.* 5 (2020) e000488.
- [21] A. Bozkurt Oflaz, B. Ekinci Köktekir, S. Okudan, Does cataract surgery simulation correlate with real-life experience?, *Turk. J. Ophthalmol.* 48 (2018) 122–126.
- [22] S. L. Cremers, A. N. Lora, Z. K. Ferrufino-Ponce, Global rating assessment of skills in intraocular surgery (GRASIS), *Ophthalmology* 112 (2005) 1655–1660.
- [23] R. Roohipoor, M. Yaseri, A. Teymourpour, C. Kloek, J. B. Miller, J. I. Loewenstein, Early performance on an eye surgery simulator predicts subsequent resident surgical performance, *Journal of Surgical Education* 74 (2017) 1105–1115. URL: <https://www.sciencedirect.com/science/article/pii/S1931720416304020>. doi:https://doi.org/10.1016/j.jsurg.2017.04.002.
- [24] A. S. S. Thomsen, P. Smith, Y. Subhi, M. I. Cour, L. Tang, G. M. Saleh, L. Konge, High correlation between performance on a virtual-reality simulator and real-life cataract surgery, *Acta Ophthalmol.* 95 (2017) 307–311.
- [25] N. PJ, Evaluation of the use of the eyesi virtual reality surgical simulator by residents and medical specialists in the argentine council of ophthalmology, *Journal of Surgery: Open Access* 6 (2020). doi:10.16966/2470-0991.213.
- [26] M. F. Jacobsen, L. Konge, D. Bach-Holm, M. la Cour, L. Holm, K. Højgaard-Olsen, H. Kjørbo, G. M. Saleh, A. S. Thomsen, Correlation of virtual reality performance with real-life cataract surgery performance, *J. Cataract Refract. Surg.* 45 (2019) 1246–1251.
- [27] L. Lucas, S. A. Schellini, A. C. Lottelli, Complications in the first 10 phacoemulsification cataract surgeries with and without prior simulator training, *Arq. Bras. Oftalmol.* 82 (2019) 289–294.

- [28] J. D. Ferris, P. H. Donachie, R. L. Johnston, B. Barnes, M. Olaitan, J. M. Sparrow, Royal college of ophthalmologists' national ophthalmology database study of cataract surgery: report 6. the impact of EyeSi virtual reality training on complications rates of cataract surgery performed by first and second year trainees, *Br. J. Ophthalmol.* 104 (2020) 324–329.
- [29] C. A. McCannel, D. C. Reed, D. R. Goldman, Ophthalmic surgery simulator training improves resident performance of capsulorhexis in the operating room, *Ophthalmology* 120 (2013) 2456–2461.
- [30] M. S. Yasar, H. Alemzadeh, Real-time context-aware detection of unsafe events in robot-assisted surgery, in: 2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), 2020, pp. 385–397. doi:[10.1109/DSN48063.2020.00054](https://doi.org/10.1109/DSN48063.2020.00054).
- [31] N. Sirajudeen, M. Boal, D. Anastasiou, J. Xu, D. Stoyanov, J. Kelly, J. W. Collins, A. Sridhar, E. Mazomenos, N. K. Francis, Deep learning prediction of error and skill in robotic prostatectomy suturing, *Surgical Endoscopy* (2024). doi:[10.1007/s00464-024-11341-5](https://doi.org/10.1007/s00464-024-11341-5), epub ahead of print.
- [32] M. B. Eppler, A. S. Sayegh, M. Maas, A. Venkat, S. Hemal, M. M. Desai, A. J. Hung, T. Grantcharov, G. E. Cacciamani, M. G. Goldenberg, Automated capture of intraoperative adverse events using artificial intelligence: A systematic review and meta-analysis, *Journal of Clinical Medicine* 12 (2023) 1687. doi:[10.3390/jcm12041687](https://doi.org/10.3390/jcm12041687).
- [33] Z. Zhong, M. Martin, M. Voit, J. Gall, J. Beyerer, A survey on deep learning techniques for action anticipation, 2023. URL: <https://arxiv.org/abs/2309.17257>. arXiv:2309.17257.
- [34] B. Fernando, S. Herath, Anticipating human actions by correlating past with the future with jaccard similarity measures, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 13219–13228. doi:[10.1109/CVPR46437.2021.01302](https://doi.org/10.1109/CVPR46437.2021.01302).
- [35] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model and the kinetics dataset, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 4724–4733. doi:[10.1109/CVPR.2017.502](https://doi.org/10.1109/CVPR.2017.502).
- [36] Z. Zhong, D. Schneider, M. Voit, R. Stiefelhagen, J. Beyerer, Anticipative feature fusion transformer for multi-modal action anticipation, 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) (2022) 6057–6066. URL: <https://api.semanticscholar.org/CorpusID:253097974>.
- [37] M. Boels, Y. Liu, P. Dasgupta, A. Granados, S. Ourselin, Supra: Surgical phase recognition and anticipation for intra-operative planning, 2024. URL: <https://arxiv.org/abs/2403.06200>. arXiv:2403.06200.
- [38] R. C. Gonzalez, R. E. Woods, *Digital Image Processing* (3rd Edition), Prentice-Hall, Inc., USA, 2006.
- [39] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by back-propagation, 2015. URL: <https://arxiv.org/abs/1409.7495>. arXiv:1409.7495.
- [40] D. Kim, Y.-H. Tsai, B. Zhuang, X. Yu, S. Sclaroff, K. Saenko, M. Chandraker, Learning cross-modal contrastive features for video domain adaptation, 2021. URL: <https://arxiv.org/abs/2108.11974>. arXiv:2108.11974.
- [41] J. Yang, H. Chen, Y. Xu, Z. Shi, R. Luo, L. Xie, R. Su, Domain adaptation for degraded remote scene classification, in: 2020 16th International Conference on Control, Automation, Robotics and Vision (ICARCV), 2020, pp. 111–117. doi:[10.1109/ICARCV50220.2020.9305483](https://doi.org/10.1109/ICARCV50220.2020.9305483).
- [42] Z. Yao, A. Gholami, D. Arfeen, R. Liaw, J. Gonzalez, K. Keutzer, M. Mahoney, Large batch size training of neural networks with adversarial training and second-order information, 2020. URL: <https://arxiv.org/abs/1810.01021>. arXiv:1810.01021.
- [43] B. Sun, K. Saenko, Deep coral: Correlation alignment for deep domain adaptation, 2016. arXiv:1607.01719.
- [44] A. Gretton, K. Borgwardt, M. J. Rasch, B. Scholkopf, A. J. Smola, A kernel method for the two-sample problem, 2008. arXiv:0805.2368.
- [45] I. Misra, C. L. Zitnick, M. Hebert, Shuffle and learn: Unsupervised learning using temporal order verification, 2016. URL: <https://arxiv.org/abs/1603.08561>. arXiv:1603.08561.
- [46] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, Wavenet: A generative model for raw audio, 2016. URL: <https://arxiv.org/abs/1609.03499>. arXiv:1609.03499.
- [47] C. Feichtenhofer, X3d: Expanding architectures for efficient video recognition, 2020. URL: <https://arxiv.org/abs/2004.04730>.
- [48] E. DeLong, D. M. DeLong, D. E. Clarke-Pearson, Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach, *Biometrics* 44 (1988) 837–845.