# IndustReal: A Dataset for Procedure Step Recognition Handling Execution Errors in Egocentric Videos in an Industrial-Like Setting

Tim J. Schoonbeek[1], Tim Houben[1], Hans Onvlee[2], Peter H.N. de With[1], Fons van der Sommen[1]

[1]Eindhoven University of Technology, Netherlands [2]ASML Research, Netherlands

t.j.schoonbeek@tue.nl

## Abstract

*Although action recognition for procedural tasks has received notable attention, it has a fundamental flaw in that no measure of success for actions is provided. This limits the applicability of such systems especially within the industrial domain, since the outcome of procedural actions is often significantly more important than the mere execution. To address this limitation, we define the novel task of* procedure step recognition (PSR), *focusing on recognizing the correct completion and order of procedural steps. Alongside the new task, we also present the multi-modal* IndustReal *dataset. Unlike currently available datasets,* IndustReal *contains procedural errors (such as omissions) as well as execution errors. A significant part of these errors are exclusively present in the validation and test sets, making* IndustReal *suitable to evaluate robustness of algorithms to new, unseen mistakes. Additionally, to encourage reproducibility and allow for scalable approaches trained on synthetic data, the 3D models of all parts are publicly available. Annotations and benchmark performance are provided for action recognition and assembly state detection, as well as the new PSR task.* IndustReal, *along with the code and model weights, is available at:*
https://github.com/TimSchoonbeek/IndustReal.

## 1. Introduction

Imagine an engineer who has just finished a service action on an internal combustion engine, only to discover that a step was missed early in the procedure. The service engineer has to undo most of the work completed after the mistake to rectify the forgotten step. The correct servicing of an engine is an example of a procedure, i.e., a given set of instructions, describing the procedural actions required to complete a task. An algorithm would become of significant value if it can understand procedural actions by automatically recognizing and tracking steps during the execution of a procedure. A system containing such an algorithm can warn users about potential mistakes or forgotten steps [3],

and summarize the execution of a procedure, thereby eliminating the need for manual logbook keeping [32]. Understanding procedures is not only highly relevant for industrial tasks, but to a broad range of procedural tasks, such as tracking the stages of surgeries to enhance post-surgical assessment and optimizing procedure workflow [15, 42].

Automated understanding of procedures is a difficult task in computer vision for various reasons. Firstly, there is often a limited visual difference between subsequent steps. For instance, determining whether a screw is correctly mounted into a specific component requires a fine-grained visual understanding. Additionally, there can be a high degree of symmetry and a low degree of textural variation for objects within procedures, especially for industrial actions [17]. Secondly, it is not feasible to collect a vast amount of data for many procedures, as they are often infrequently performed and rather specialized. Finally, procedures can typically be completed correctly via several possible execution orders. Therefore, it is frequently not sufficient to merely look whether a single, pre-defined step is completed correctly.

Existing approaches to procedural understanding can generally be divided into two groups, one performing action recognition (AR) [30, 33] and the other assembly state detection (ASD) [14, 27, 36]. AR approaches aim to recognize only which actions are being performed, rather than which actions are actually completed. This is a crucial difference, and it is more valuable to know whether a step has been actually completed correctly, than to know only if an operator spend time on that step. ASD relies on object detection algorithms, detecting the actual phase during an object's assembly within a procedure. However, the number of possible states in ASD explodes with the number of parts in a procedure, limiting the complexity in existing literature to five or six parts only [14, 27, 36]. Furthermore, both AR and ASD implementations do not explicitly leverage procedural knowledge, *e.g.* which actions are to be expected after the observation of a preceding action.

To address the aforementioned limitations, this work formally defines the novel task of *procedure step recog-*
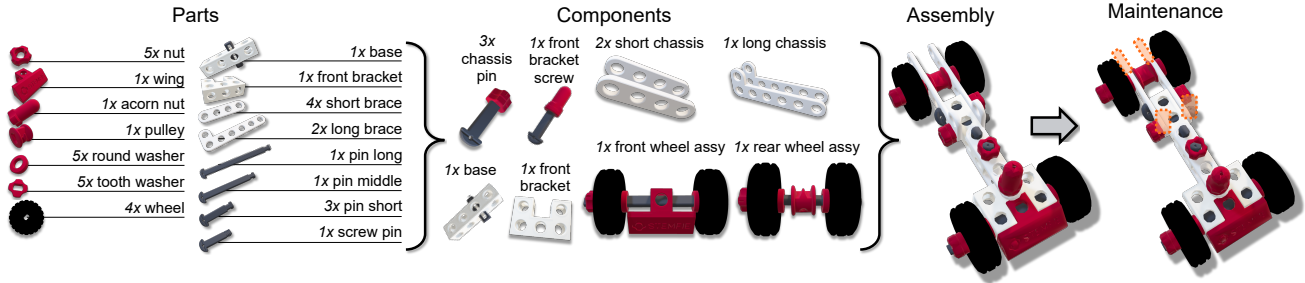
Figure 1. Overview of the construction-toy car's (3D printed) parts, components, and models used in IndustReal. The maintenance task represents a component upgrade and consists of replacing the long braces with short braces on the rear chassis.

*nition* (PSR) and introduces IndustReal, a publicly available dataset towards solving this task. It is an ego-centric, multi-modal dataset where 27 participants are challenged to perform assembly and maintenance procedures on a construction-toy car, based on STEMFIE [22], demonstrated in Fig. 1. The videos are annotated for action recognition (AR), assembly state detection (ASD), and procedure step recognition. The placement of IndustReal within literature is outlined in Tab. 1. The IndustReal dataset features four novel aspects:

- *Variety of execution errors.* IndustReal features 38 errors, of which 14 are exclusive to the validation and test sets. Whilst some datasets already include procedural errors (*e.g.* omissions) [28, 33], IndustReal is the first to also include execution errors (*e.g.* wrong type of nut used).

- *Subgoal oriented execution.* Other datasets are either "free-style" assemblies [33] or contain a strict, step-by-step execution order [30]. IndustReal combines these execution types with a subgoal-oriented assembly style, where participants are given flexibility to determine the execution order between subgoals. This approach more closely resembles industrial procedures, since it maintains a hierarchy in procedure execution whilst allowing for flexibility where possible. IndustReal contains 48 different execution orders.

- *Open-source geometries.* Scalability is an important factor for many industrial tasks, where simultaneously the technical drawings are often available. Therefore, 3D models for all parts are published, to stimulate use of synthetic data in procedural action understanding, *e.g.* by sim2real domain adaptation or generalization.

- *3D printed parts.* To ensure reproducibility, future availability of the model, and growth via community effort, all parts are 3D printed and open source [22].

In summary, the contribution of this work is two-fold: we present the IndustReal dataset, define the task of procedure step recognition and provide a benchmark towards this task.

## 2. Related Work

### 2.1. Work on procedural action understanding

In (procedural) action recognition tasks, the objective is to classify a video clip into a set of activities expected within a certain procedure [14, 27, 36, 38]. Examples of such procedures are cooking [4, 10, 41], instruction videos [26, 34], and assembly tasks [16, 18]. Recently, focus has been directed towards recognizing procedural actions in industrial-like settings [2, 6, 28, 30, 33].

Whilst recognizing actions in procedural videos is certainly of some interest in an industrial setting, we argue that it is more valuable to know what was actually completed, rather than only performed (and therefore potentially not finished). Therefore, the PSR task is proposed to recognize the correct completion of steps and to track the order in which those steps were completed.

Additionally, the comparable datasets outlined in Tab. 1 generally have restricted sizes, since it is challenging to record sufficient video data of people executing (nearly) the same procedural task. This is especially relevant to industrial applications, where a significant portion of the tasks is performed rather infrequently. However, procedural knowledge, in the form of work instructions or technical drawings, is generally available for those tasks. Whilst such procedural knowledge is not explicitly leveraged by the aforementioned approaches, we advocate to explicitly use this knowledge in PSR to constrain the number of possible actions expected at any given moment.

### 2.2. Literature on assembly state detection

Assembly state detection is a sub-task of object detection, where the objective is to recognize and locate the specific state of an object during an assembly procedure [14, 27, 35, 36, 39, 40]. ASD is significantly more challenging than object detection, since the visual variability between two states is often small. For instance, detecting the object "Car" is significantly less complicated than detecting the state "Car without window". Su *et al.* [36] simul-

| Dataset | Year | Ego | AR | Tasks ASD | PSR | Flex. | PEs | EEs | Parts | 3DM | Seqs. | Dur. | Participants |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IKEA ASM [6] | 2021 | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | 7 | ✗ | 381 | 35.3h | 48 |
| MECCANO [29, 30] | 2021 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓* | ✗ | 49 | ✗ | 20 | 6.9h | 20 |
| Assembly101 [33] | 2022 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓* | 15 | ✗ | 1.01K | 167.0h | 53 |
| BRIO-TA [28] | 2022 | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | 10 | ✗ | 75 | 2.9h | 15 |
| HA4M [9] | 2022 | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | 17 | ✓ | 217 | 5.9h | 41 |
| ATTACH [2] | 2023 | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | 26 | ✗ | 378 | 17.2h | 42 |
| IndustReal (ours) | 2024 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 36 | ✓ | 84 | 5.8h | 27 |

Table 1. Comparison of industrial-like procedural understanding datasets. AR: action recognition, ASD: assembly state detection, PSR: procedure step recognition, flex.: more than one single execution order for the task, PEs: procedural errors (omission, execution order), EEs: execution errors (component installed incorrectly), 3DM: 3D models made publicly available, *rare and not explicitly annotated.

taneously detect the state and the pose of a coffee machine during its assembly. The authors demonstrate a good performance on the task, but the different assembly states consist of visually distinctive objects. Lui *et al.* [27] demonstrate an attention mechanism for their convolutional neural network (CNN), detecting object states even if they are visually similar to each other. Both works demonstrate promising advances on this task. The approaches are trained predominantly on synthetic data, which can be readily generated if 3D models of each assembly state are already designed.

Nevertheless, the aforementioned approaches have two important limitations. Firstly, the algorithms must be trained on a dataset of images from each object state. This forces the models to learn a low-dimensional representation for each state. Since the procedures that the authors selected consist of only 5 parts and, at most, 6 distinctive object states, training a model in such a manner is feasible. However, it remains unclear whether this approach scales to procedures with higher complexity and more object states. Secondly, no procedural information is leveraged by either approach. Therefore, the models expect for example, the initial state in an assembly equally as much as the very last state. For tasks with sufficient visual distinction between states, this is not a critical limitation. However, for objects that appear identical from certain viewpoints for several different states, procedural information could be used to determine the most likely state. Finally, none of the above-mentioned approaches release their (test) data.

### 2.3. Published comparable datasets

The placement of IndustReal within publicly available datasets, recorded in industrial-like settings is shown in Tab. 1. The MECCANO dataset [29, 30] is most relevant to IndustReal in terms of procedure, task complexity, and dataset size. Notable differences between these two datasets are that in MECCANO users follow strict, step-by-step instructions, limiting the variety in execution order. Secondly, although MECCANO contains some procedure errors, they

are nearly always corrected in the subsequent step. Therefore, later states do not have prior errors in them, whilst such cases are certainly of interest and are frequently encountered in the industrial domain. Finally, we argue that the availability of 3D models of all parts represented in the dataset is crucial for industrial applications, given the wide availability and usage of technical drawings in that domain. Unfortunately, the MECCANO dataset uses a IP-protected construction set, prohibiting the publication of such CAD models. BRIO-TA [28] and the large-scale Assembly101 [33] datasets both contain procedural mistakes, such as omissions and incorrect execution order, but do not contain (labeled) execution errors. Furthermore, none of the mentioned datasets provide labels for assembly state detection. Lastly, to ensure future availability of the models, all parts in IndustReal are 3D printed. An added benefit of 3D printing is that it allows researchers to print in different scale, colors, or materials, to test the limits of their algorithms.

## 3. Procedure Step Recognition

The previous section has identified a gap in the related work between AR and ASD. By formalizing the task of procedure step recognition (PSR), along with an evaluation scheme, we encourage researchers to develop methods to automatically recognize the completion of steps, rather than the (partial) execution of activities. Additionally, PSR systems should explicitly leverage procedural knowledge and allow a flexible execution order for procedural tasks, when the procedure allows it.

### 3.1. Task definition

The objective of PSR is to extract an estimate of all procedure steps correctly performed by a person up to time $t$, based on sensory inputs $X_t = (x_t, x_{t-1}, \ldots, x_{t-h})$ and a descriptive set of the procedural actions to be performed $\mathcal{P} = \{a_0, a_1, \ldots, a_n\}$. Here, $h$ is the observation horizon and $n + 1$ the total number of actions $a_i \in \mathcal{P}$ covered in

the procedure. The predicted completed procedure steps $\hat{y}_t$ at time $t$, given some computational model $\mathcal{F}$, are defined such that

$$\hat{y}_t = \mathcal{F}(X_t, \mathcal{P}). \tag{1}$$

Here, a "predicted" (recognized) procedure step does not refer to the prediction of a future step, but rather the prediction by a model for a correctly completed step, based on the given inputs. Crucially, this definition allows for real-time operation, since contrary to existing tasks, PSR does not require a full recording of the procedure as input [23, 26].

Each unique action $a_i \in \mathcal{P}$ contains information regarding this specific action, and can be different for varying approaches to PSR. For instance, $a_i$ can contain a step description, a template image of what the completed step is supposed to look like, or relative positions between components. Sensory inputs $X_t$ may comprise of camera images, depth maps, or even user interaction with the system, during the execution of a procedure. Note that with the given definition, a PSR system can be an ensemble of various computer vision algorithms, *e.g.* object detection to determine the assembly state and action recognition to determine the activities.

The predicted procedure steps form an ordered list of elements, describing the completed steps and is defined as

$$\hat{y}_t = (\hat{s}_{\sigma(0)}, \hat{s}_{\sigma(1)}, \ldots, \hat{s}_{\sigma(m)}), \tag{2}$$

where step $\hat{s}_{\sigma(i)}$ is the predicted completion of the action $a_i \in \mathcal{P}$, at prediction time $\hat{t}_{\sigma(i)}$, having a prediction confidence $c_{\sigma(i)}$, with $m + 1$ the total number of recognized procedure steps. The function $\sigma$ maps a completed step $\hat{s}_j \in \hat{y}_t$ to the corresponding action $a_i \in \mathcal{P}$. Therefore, the first predicted step $\hat{s}_{\sigma(0)}$ is not necessarily the completion of the action $a_0 \in \mathcal{P}$, but rather the first recognized step that a person completed.

The ground-truth execution order $y_t$ at time $t$ describes the order in which the steps are actually completed, which can differ from the prescribed order in $\mathcal{P}$, and is defined as

$$y_t = (s_{\rho(0)}, s_{\rho(1)}, \ldots, s_{\rho(k)}), \tag{3}$$

where $s_{\rho(i)}$ is the completion of the action $a_i \in \mathcal{P}$ at time $t_{\rho(i)}$. The value $k + 1$ is the total number of actions completed and $\rho$ a function that maps the completed steps from $y_t$ to the steps described in $\mathcal{P}$. If the order of the steps predicted in $\hat{y}_t$ equals that of $y_t$, it follows that $\sigma = \rho$, signifying a perfect execution order prediction.

## 3.2. Evaluation metrics

To quantify the performance of a PSR system, three evaluation metrics are proposed. These metrics focus on predicting the procedure steps in the correct order, the number of false predictions, and the timeliness of the predictions.

**Metric 1: procedure order similarity.** We propose to measure the quality of a predicted sequence order for an entire recording, $\hat{y}$ (e.g., 'ACB'), by comparing it with a similarity measure with respect to the ground-truth $y$ (e.g., 'ABC'). This is approached as a string similarity problem, a common problem in spelling error detection [5, 31], since words consist of a sequence of characters, where order and type of character matters. In temporal action segmentation, the Levenshtein (Lev) distance, normalized over the length of the ground-truth sequence, is commonly used [23]. We propose two changes to this metric, by (1) eliminating substitution from the edit distance, preventing the metric from favouring models with many false positives, and (2) using the Damerau-Levenshtein (DamLev) [11], rather than the Lev distance because it penalizes transpositions less, since it is intuitive to penalize "A<u>CB</u>" less compared to "<u>CA</u>B".

In contrast to Lea *et al.* [23], we propose to normalize the edit distance with respect to the length of the ground truth, rather than the length of either the ground truth or the prediction, depending on which is longer. This prevents models with many false positives from being normalized favourably. Finally, the normalized edit distance is subtracted from the unity value, resulting in a similarity metric, rather than a distance metric. Thus, the procedure order similarity (POS) between $y$ and $\hat{y}$ is defined as

$$\text{POS} = 1 - \min\left(\frac{\text{DamLev}(y, \hat{y})}{|y|}, 1\right), \tag{4}$$

where $\text{DamLev}(\cdot)$ is a weighted DamLev edit distance function. Further clarification on POS may be found in [1].

**Metric 2: $F_1$ score.** A false positive is defined as a procedure step $\hat{s}_{\sigma(j)}$ that is predicted prior to the actual completion of action $a_i$, or if $a_i$ is not at all completed, hence

$$(\hat{t}_{\sigma(j)} < t_{\rho(i)}) \vee (a_i \notin y). \tag{5}$$

A false negative is defined as a step $s_{\rho(j)}$ for a corresponding action $a_i$, that has indeed been completed, but is not represented in $\hat{y}$, so that

$$(a_i \in y) \wedge (a_i \notin \hat{y}). \tag{6}$$

Finally, a true positive is defined as the prediction of a procedure step $\hat{s}_{\sigma(j)}$, that is observed at or after the actual completion of $a_i$, such that

$$(\hat{t}_{\sigma(j)} \geq t_{\rho(i)}) \wedge (a_i \in y). \tag{7}$$

PSR systems can explicitly assume that actions have been completed based on procedural information, without relying exclusively on sensory recognition. Therefore, a distinction can be made between $F_1$ score at the recognition and system level. Since this definition does not contain any
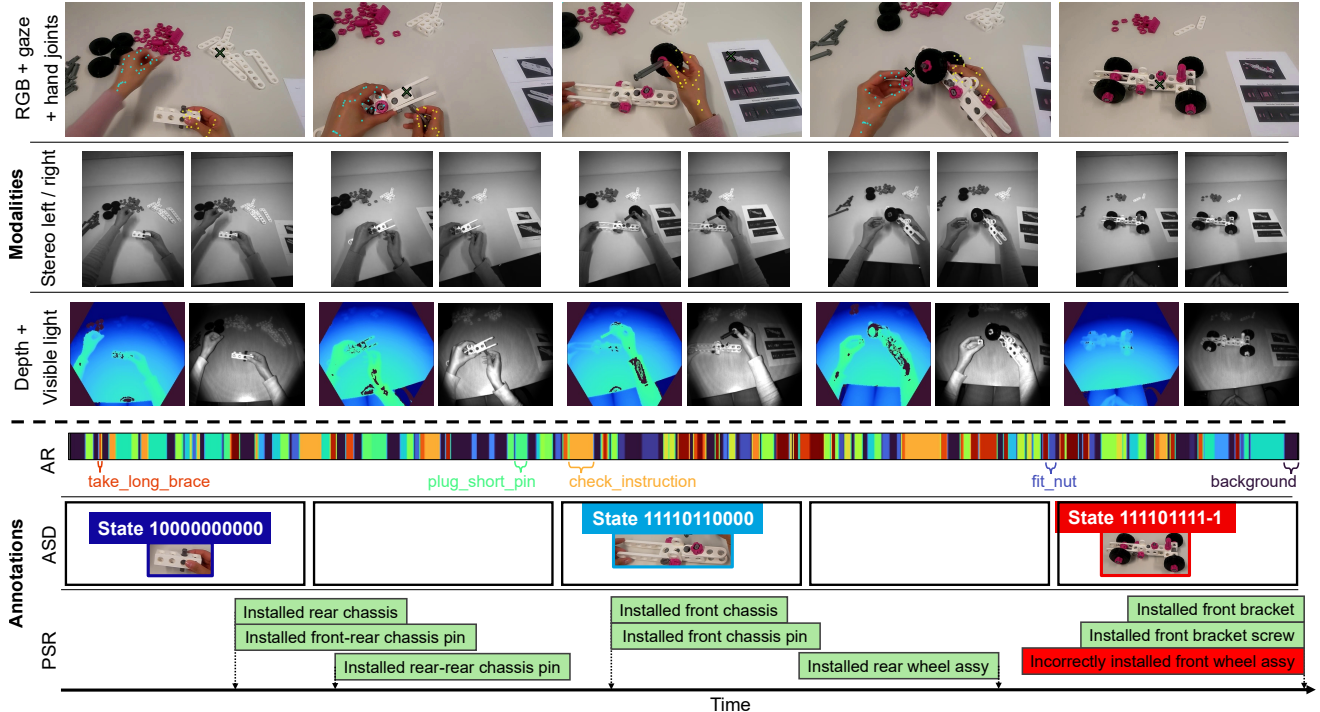
Figure 2. Samples from a clip in the IndustReal dataset, demonstrating the modalities and annotations for all three tasks. Gaze is indicated by the cross, detected hand joints by the dots. AR: action recognition, ASD: assembly state detection, PSR: procedure step recognition.

time restriction on true positives, procedure steps that are recognized long after the step completion are not penalized. Therefore, comparing only the procedure order similarity and $F_1$ score lacks a temporal component.

**Metric 3: average delay.** To complement the aforementioned metrics with a temporal component, the average delay $\tau$ is introduced, quantifying the time between the ground-truth completion and corresponding recognition of a step. False negatives, defined by Equation 6, have an undefined delay, and similarly, false positives (defined in Equation 5) have either an undefined delay, or a negative delay. Therefore, FPs and FNs are excluded from $\tau$, defines as

$$\tau = \frac{1}{h} \sum_{i=0}^{h-1} (\hat{t}_{\sigma(i)} - t_{\rho(i)}), \tag{8}$$

where $h$ is the number of total TPs in $y$. Because FPs and FNs are discarded in the delay, only the combination of the three proposed metrics provides a valuable insight into the performance of a system towards solving PSR.

## 4. IndustReal Dataset

### 4.1. Construction-set car procedures

Figure 1 demonstrates the 36 part models used in the IndustReal dataset, based on the STEMFIE construction-

set toy car [22]. The model has significant complexity, consisting of multiple types of washers, pins, and braces. Some components require screwing, others need tightening, and participants frequently have to use both hands. Two procedures are defined, an assembly task, where the car has to be build from scratch, and a maintenance task, where the participants have to replace part of the rear chassis of the toy car. Printed instructions are provided in a subgoal-oriented manner, meaning that the participant builds towards subgoals rather than executing strict, step-by-step instructions, or "free-style" building towards a final assembly. Unlike related datasets, participants are allowed to create subassemblies, as commonly encountered in industrial settings.

All of the parts are 3D printed on an Ultimaker S5 at 200% scale, layer height of 0.3 mm, 15% infill, a print speed of 50 mm/s, and PLA filament and the colors white, silver metallic, magenta and black. All parts, as well as the final models, are published together with the dataset, since part geometries are commonly available in industrial settings, *e.g.* CAD models.

### 4.2. Recording and setting

The HoloLens 2 (HL2) [37] mixed reality headset is used as recording rig for the dataset. The front-facing RGB camera records at a resolution of 1280×720 pixels and the stereo cameras provide images at 480×620 pixels. The long-throw depth and IR sensors record at a resolution of

320×288 pixels and contain normalized values. Next to the images, we also record gaze, hand, and head-pose tracking, provided by the HL2 algorithms. All sensors, visualized in Fig. 2, are sampled at 10 fps, except from the depth and IR sensors, as these are limited to 5 fps in hardware. The data are sent from the HL2 to a server in real time, using the HL2SS library [19].

The dataset is recorded in a setting with consistent background and lighting conditions. More specifically, participants are asked to perform the procedures on a white desk placed against a white background. Such lack of variety in background and lighting conditions is assumed to be consistent with industrial settings. Redundant washers, nuts, screws, and pins are placed amongst the required parts to reduce bias towards detecting unused parts. Prior to each recording, all parts are sorted into heaps according to color, and then placed randomly on the desk.

## 4.3. Participants and protocol

In total, 27 participants were recruited. Each participant signed a consent form and the experiment was approved by the institution's Ethical Review Board. Each participant was asked to perform the assembly procedure once, without recording, to familiarize themselves with the construction-set car. During this practice assembly, feedback was provided to the participants when required. Subsequently, the HL2 was fitted to the participant and the gaze tracking was calibrated. Then, each participant was assigned a correct assembly instruction, a correct maintenance instruction, and one or two instructions with errors, to ensure that a variety of mistakes were introduced. As anticipated, the participants exhibited errors even when they were provided with the correct instructions, and these mistakes were annotated as well. The recorded errors in the dataset vary from minor, difficult-to-observe mistakes (*e.g.* installing the wheels without washers), to large mistakes, such as forgetting a wheel. To minimize "learners bias" of the participants in the dataset, each participant was given the instructions in a random order. Participants were asked to remove their hands from the assembly after completing a step, to ensure a minimally occluded view on the assembly state.

## 4.4. Annotation

Since IndustReal is intended to be closely representative of industrial use cases, significant importance is given to evaluation and robustness to unseen, out-of-distribution errors. Additionally, a large real-world test set is desired when training exclusively on synthetic data. Therefore, a dataset split of 12/5/10 (No. of train/val/test participants) is chosen, which is heavily focused on the test set. Furthermore, numerous errors are exclusively present in the validation or test sets. Note that the split is made on participants, rather than videos, to ensure sufficient variation in viewpoint, ex-

ecution, and head movements between the three sets.

Although IndustReal is specifically presented to address PSR, annotations are also provided for AR and ASD, enabling researchers to use IndustReal for various tasks, or to combine tasks for better PSR performance in future work.

### 4.4.1 Action recognition

AR labels consist of the frame at which the action starts, ends, and the combination of a verb and noun. Given the resemblance between MECCANO [30] and the construction-toy car used in IndustReal, we adopt the same verbs as utilized in MECCANO: *take*, *put*, *align*, *plug*, *pull*, *screw*, *unscrew*, *tighten*, *loosen*, *fit*, *check*, and *browse*. Using these verbs, the component names described in Fig. 1, and the additional nouns *objects*, *partial model*, and *instruction*, we have annotated 75 fine-grained action classes and a total of 9,273 instances. The average action lasts 1.9±1.4 seconds. A long-tail distribution is observed, with 80% of the data containing 29.3% of the actions, and provided together with more statistical details in [1]. Recordings have an average of 110±38 actions per video, with 134±32 actions per assembly and 79±13 actions maintenance video. Participants frequently perform actions simultaneously, resulting in 24.2% of all instances having an overlap with at least one other action.

### 4.4.2 Assembly state detection

To the best of our knowledge, IndustReal is the first publicly available dataset with ASD annotations. Because the participants have flexibility in execution order, all labeled states must be explicitly defined. We label the assembly states with integers, where a "1" indicates that a component at that index has been correctly installed and a "0" indicates that it has not (yet) been correctly installed. Additionally, we assign "-1" to components that have been incorrectly installed. Rather than labeling each assembly part individually with such a code, we divide the toy car into 11 components (in order): *base*, *front chassis*, *front chassis pin*, *rear chassis*, *short-rear chassis*, *front-rear chassis pin*, *rear-rear chassis pin*, *front bracket*, *front bracket screw*, *front wheel assy*, and *rear wheel assy* (see Fig. 1). For example, the assembly state *11100000000* consists of a correctly installed base, front chassis and front chassis pin.

Bounding boxes and labels are provided for all 22 labeled (defined) states in IndustReal, as well as 27 different error states. Intermediate states, which occur during the assembly of components, are not labeled, as outlined in Fig. 2. Whilst those states could be annotated as *partial model*, such labels empirically appear to hold little value. Intermediate states differ from erroneous states in that participants in intermediate states are actively progressing towards completing a state, whereas participants in erroneous

| Model | Modalities | Top-1 acc. [%] | Top-5 acc. [%] |
|---|---|---|---|
| SlowFast [13]* | RGB | 57.83 | 82.87 |
| SlowFast [13]† | RGB | 60.39 | 85.21 |
| MViTv2 [24]* | RGB | 62.43 | 85.62 |
| MViTv2 [24]† | RGB | 65.25 | 87.93 |
| SlowFast [13]† | RGB, VL, stereo | 62.34 | 85.97 |
| MViTv2 [24]† | RGB, VL, stereo | **66.45** | **88.43** |

Table 2. AR benchmark on IndustReal. *MECCANO [30] pre-trained, †Kinetics [21] pre-trained, VL: visible light.

| Pre-trained | Fine-tuned | mAP (b-boxed) | mAP (entire videos) |
|---|---|---|---|
| COCO | Synthetic | 0.573 | 0.341 |
| COCO | IndustReal | 0.753 | 0.553 |
| Synthetic | IndustReal | 0.779 | 0.575 |
| COCO | IndustReal + synthetic | **0.838** | **0.641** |

Table 3. ASD performance benchmark on IndustReal using YOLOv8-m [20] for various training schemes.

states are not trying to further complete that step. In total, 26.9K video frames (13% of total) are annotated for ASD, of which 3,569 frames show error states.

### 4.4.3 Procedure step recognition

The PSR labels consist of the frame at which a step completion occurs and the new assembly state, as defined in the previous section. The difference between two assembly states can directly be used to determine which procedure steps are completed. A procedure step is defined as *completed* when the component relating to that step is correctly installed, which includes actions such as the tightening of a nut. Although the explicit detection of incorrect step execution is not included in the PSR task definition (Sec. 3.1), annotations of incorrect assembly states are included for qualitative analysis and future work. Therefore, two sets of PSR labels are provided, one with only correctly executed procedure steps (*e.g.* "Installed front chassis"), and one that also includes incorrectly completed steps, such as "Incorrectly installed front wing". An example of a PSR-annotated sequence is outlined in Fig. 2. In total, IndustReal consists of 724 correct procedure step completions (8.6±1.2 correct completions per recording) and 38 incorrect step completions. In total, 35 videos (42%) contain a missing or incorrectly completed procedure step. IndustReal contains 22 different correctly completed procedure execution orders, plus an additional 26 different execution orders containing error states.

## 5. Benchmark Experiments

This section outlines the benchmark performance of state-of-the-art approaches to AR and ASD on IndustReal. Furthermore, a PSR benchmark implementation is outlined and evaluated, providing a baseline performance for more sophisticated approaches towards this task.

### 5.1. Action recognition benchmark

**Definition.** Given a video segment $X_i = [x_{ts_i}, x_{te_i}]$ and a set of action classes $C_a = \{c_0, c_1, ..., c_n\}$, the objective of action recognition is to classify segment $X_i$ to the correct class $c_i \in C_a$ [30]. Here, $x_{ts_i}$ and $x_{te_i}$ indicate the start and end frame for action $c_i \in C_a$, respectively.

**Benchmark.** For this task, the SlowFast [13] CNN and MViTv2-S [24] transformer are chosen to benchmark the performance. Each model is trained on IndustReal after Kinetics [21] pre-training. Additionally, since the MEC-CANO dataset for action recognition is closely related to the IndustReal dataset, we also report baselines pre-trained on MECCANO. Finally, both networks are trained on depth, visible light (VL), and stereo images, and combined to create an ensemble of models trained on various modalities.

Top-1 and top-5 accuracy are reported in Tab. 2 for the aforementioned experiments. Pre-training on the MEC-CANO dataset does not provide a performance benefit. The MViTv2 transformer outperforms the SlowFast architecture. The best performance is observed for the MViTv2 ensemble of the modalities RGB, VL, and stereo images. As motivated in [1], depth is excluded from this ensemble. Notably, for both architectures, each individual modality is outperformed by the ensemble, indicating that each modality contains some complementary form of relevant information.

### 5.2. Assembly state detection benchmark

**Definition.** Given a video frame $X_i$ and a set of assembly states $Z_a = \{z_0, z_1, ..., z_n\}$, the objective of ASD is to detect the bounding box and assembly state $z_i \in Z_a$ for the sample $X_i$. The states $Z_a$ are defined in Sec. 4.4.2.

**Benchmark.** For this benchmark, the state-of-the-art object detection network YOLOv8-m [20] is employed. Since geometries of all parts are published with IndustReal to stimulate synthetic learning, we provide four training schemes. First, we train the model, pre-trained on COCO [25], exclusively on a synthetically generated dataset. To generate the synthetic data, Unity Perception [7] is used to generate 100K training samples, each containing one assembly state in $Z_a$. Secondly, we train the model (pre-trained on COCO) directly on IndustReal. Then, we pre-train on the synthetic dataset, after which we fine-tune on the IndustReal dataset. Finally, the synthetic and real-world datasets are combined for the last baseline.

|      | All recordings | | | Recordings with errors | | |
|------|------|-------|------------|------|-------|------------|
|      | POS  | $F_1$ | $\tau$ [s] | POS  | $F_1$ | $\tau$ [s] |
| B1   | 0.570 | 0.779 | **14.9** | 0.480 | 0.698 | **14.4** |
| B1-S | 0.014 | 0.206 | 36.9 | 0.000 | 0.174 | 48.4 |
| B2   | 0.731 | 0.860 | 22.3 | 0.636 | 0.784 | 20.2 |
| B2-S | 0.240 | 0.573 | 44.4 | 0.107 | 0.516 | 60.5 |
| B3   | **0.797** | **0.883** | 22.4 | **0.731** | **0.816** | 20.4 |
| B3-S | 0.597 | 0.734 | 49.5 | 0.571 | 0.731 | 71.4 |

Table 4. PSR performance benchmark on IndustReal with ASD backbone, expressed by procedure order similarity (POS), $F_1$ score, and average delay $\tau$. B1-3 denote the three baseline models and S indicates exclusive training on synthetic data.



(a) ASD predict: base, rear-rear pin and rear chassis correctly installed.



(b) ASD predict: entire assembly procedure correctly installed.

Figure 3. Visualization of ASD and PSR results for two cases. Bounding boxes and captions indicate the ASD predictions, text boxes indicate PSR predictions (B3). Both ASD predictions are false positives. (a) Front-rear pin orientation is incorrect, whilst the ASD predicts the absence of this pin. (b) Nut instead of screw is incorrectly used for the front brace, which is not recognized.
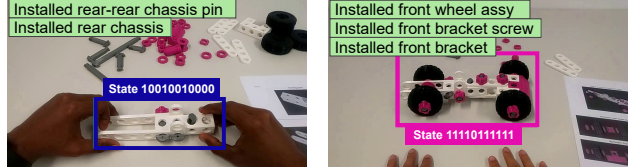
The results are quantified using the mAP metric and reported on frames of the IndustReal test set containing ground-truth bounding-box annotations, as well as the entire test set. As outlined in Tab. 3, combining the synthetic and real-world images results in the highest performance. A significant performance drop of 27% is observed when evaluating entire videos, rather than only frames with a ground-truth annotation. The decreased performance is caused by false positive predictions on states with fine-grained visual differences, predominantly error states and directly prior to the completion of a procedural step. Notably, the best performing model has a false positive rate of 65% and average precision (AP) of 0.23 for assembly states containing an error. Two such error states, and their corresponding ASD predictions, are visualized in Fig. 3, and further qualitative and quantitative analysis is provided in [1].

### 5.3. Procedure step recognition benchmark

**Definition.** The task definition is given in Sec. 3.1.
**Benchmark.** For the PSR benchmark, we provide three baseline implementations relying on the ASD model output outlined in the previous section. The first baseline B1 determines for each change in detected assembly state which corresponding steps must have been completed since the last detected state, and assumes these actions to be executed correctly. The second baseline B2 accumulates the confidence for each predicted step completion (obtained from B1) over time, until a threshold $T$ is reached, upon which an action is deemed correctly completed. Baseline 3 uses the same confidence aggregation as B2, but limits the number of possible step completions to those expected in the correct execution of the given procedure. Further clarification on the baselines may be found in [1]. For each baseline, the best performing ASD model from Tab. 3 is used. Additionally, each baseline is evaluated with the ASD approach trained exclusively on synthetic data. The detection pipeline (ASD + PSR) yields real-time operation at 178 fps on a v100 GPU.

The POS, $F_1$, and $\tau$ metrics for all baselines are outlined in Tab. 4. The best performing approach achieves relatively high POS and $F_1$ scores (0.797 and 0.883 on all recordings, respectively), due to under-representation of execution errors compared to correctly executed steps. On recordings with errors, all baselines show a distinctive decrease in performance, as shown in Fig. 3, thereby highlighting the need for approaches that better handle (out-of-distribution) errors. Although the performance of B3-S trained entirely on synthetic data is significantly lower than its real-world trained counterpart, it notably outperforms B2-S. This indicates that using procedural information to restrict the set of possible step completions can considerably improve the model performance.

## 6. Conclusion

This paper proposes the novel task of *procedure step recognition* (PSR). This task bridges the gap that currently existing *action recognition* and *assembly state detection* tasks leave in procedural activity understanding, by explicitly leveraging procedural information and recognizing correctly completed steps. This paves the way for the development of upcoming, increasingly powerful procedure-assistive systems.

Along with the PSR task, the IndustReal dataset is presented. IndustReal differs from existing datasets in the wide variety of procedural and execution errors, subgoal-oriented procedure execution, and 3D printed parts to ensure future reproducibility. Additionally, the geometries of object parts are published to stimulate sim2real domain adaptation and generalization, as this is crucial to scalability in industrial settings. Furthermore, performance benchmarks are provided for AR, ASD, and PSR based on the proposed IndustReal dataset with various pre-training techniques.

The PSR baseline shows promising results on procedural videos where participants do not make mistakes, but fails to generalize to (out-of-distribution) execution errors. Future work will focus on improving upon this baseline performance and increase scalability, such that PSR becomes viable and indeed attractive towards more industrial use cases.

# Acknowledgment

# 7. *Supplementary Material for:*

## IndustReal: A Dataset for Procedure Step Recognition Handling Execution Errors in Egocentric Videos in an Industrial-Like Setting

## 7.1. Introduction

This supplementary material accompanies the main paper titled *IndustReal: A Dataset for Procedure Step Recognition Handling Execution Errors in Egocentric Videos in an Industrial-Like Setting*. This document offers readers a comprehensive and more in-depth understanding of the conducted research.

Section 7.2 provides additional information on the action recognition (AR), such as the distribution of class instances, implementation details, and further quantitative results. Section 7.3 provides implementation details on assembly state detection (ASD) and outlines additional qualitative as well as quantitative results. Lastly, Section 7.4 provides further clarification on the motivation for the proposed procedure order similarity (POS) evaluation metric, implementation details and pseudo code describing the procedure step recognition (PSR) baselines, and a qualitative example.

## 7.2. Action recognition

### 7.2.1 Annotations

Annotators were instructed to mark the action start when the participants initiate the action, rather than when participants touch the relevant object(s). Figure 5 demonstrates the action class instances and long-tail distribution for action recognition annotations in the IndustReal dataset. As can be seen, the *check_instruction* action class is the most common, followed by *align_objects*. The long-tail distribution demonstrates that 23 action classes (out of 73) constitute 80% of the dataset. All actions were annotated using the ELAN tool [8]. Two annotators annotated the entire dataset, and each annotator reviewed the annotations of the other.



| | | |
|---|---|---|
| 0: take_short_brace | 25: plug_partial_model | 50: fit_wheel |
| 1: align_objects | 26: plug_pin_middle | 51: check_partial_model |
| 2: take_pin_short | 27: take_pulley | 52: put_short_brace |
| 3: plug_short_pin | 28: plug_wheel | 53: fit_objects |
| 4: take_tooth_washer | 29: browse_instruction | 54: put_round_washer |
| 5: take_nut | 30: fit_short_brace | 55: fit_pulley |
| 6: tighten_nut | 31: fit_tooth_washer | 56: fit_wing_beam |
| 7: check_instruction | 32: fit_round_washer | 57: put_tooth_washer |
| 8: take_partial_model | 33: fit_long_brace | 58: pull_pin_middle |
| 9: take_long_brace | 34: fit_nut | 59: put_wing_beam |
| 10: take_screw_pin | 35: put_screw_pin | 60: put_pulley |
| 11: take_instruction | 36: put_wheel | 61: pull_screw_pin |
| 12: put_instruction | 37: check_wheel | 62: put_acorn_nut |
| 13: take_pin_long | 38: pull_wheel | 63: loosen_acorn_nut |
| 14: put_pin_long | 39: loosen_nut | 64: fit_partial_model |
| 15: take_wing_beam | 40: put_nut | 65: take_small_screw_pin |
| 16: plug_screw_pin | 41: pull_objects | 66: plug_small_screw_pin |
| 17: take_round_washer | 42: put_pin_middle | 67: put_small_screw_pin |
| 18: take_acorn_nut | 43: take_objects | 68: fit_acorn_nut |
| 19: tighten_acorn_nut | 44: put_partial_model | 69: fit_wing |
| 20: take_pin_middle | 45: put_objects | 70: pull_pin_long |
| 21: take_wheel | 46: pull_pin_short | 71: plug_objects |
| 22: plug_pin_long | 47: put_pin_short | 72: pull_small_screw_pin |
| 23: take_wing | 48: put_long_brace | 73: tighten_tooth_washer |

Figure 4. Normalized confusion matrix for action recognition (MViTv2 [24] pretrained on Kinetics [21]).

### 7.2.2 Implementation details

The AR baselines (SlowFast [13] and MViTv2 [24]) are trained using the SlowFast library [13], and some baselines are pre-trained on MECCANO [30]. For all baselines, the configuration yaml files are provided on the dataset repository, such that they can directly be used to reproduce the reported results. The baselines are trained on one Nvidia Tesla v100 GPU.

### 7.2.3 Quantitative and qualitative analysis

Figure 4 shows the normalized confusion matrix for action recognition on the IndustReal test set. Specifically, it shows the performance of the MViTv2 [24] transformer, pretrained on the Kinetics dataset [21]. It is observed that although the performance is generally adequate, the model confuses actions that are visually similar, such as *take* and *put* the *short_brace*, and *tighten* and *loosen* the *acorn_nut*.

Furthermore, the performance of both the MVit and SlowFast architectures on the RGB, depth, visible light, and stereo image modalities are outlined in 5. For both archi-

Figure 5. Visualization of the distribution of action recognition labels in the IndustReal dataset.

| Model | Modality | Top-1 acc. [%] | Top-5 acc. [%] |
|---|---|---|---|
| SlowFast [13] | RGB | **60.39** | **85.21** |
| SlowFast [13] | Depth | 43.20 | 73.98 |
| SlowFast [13] | Visible light | 53.75 | 81.48 |
| SlowFast [13] | Stereo images | 57.72 | 83.03 |
| MViTv2 [24] | RGB | **65.25** | **87.93** |
| MViTv2 [24] | Depth | 49.08 | 76.51 |
| MViTv2 [24] | Visible light | 58.59 | 83.50 |
| MViTv2 [24] | Stereo images | 58.86 | 83.55 |

Table 5. AR performance benchmark on IndustReal per modality. All models are pre-trained on Kinetics [21].

tectures, RGB outperforms the other modalities, likely due to the extensive pre-training on the (RGB) Kinetics [21] dataset. Due to hardware limitations on the current HL2 operating system, it is not possible to access the short-throw depth and RGB camera simultaneously, hence the short-throw depth data are not provided. However, the stereo images can be used to generate high-resolution depth maps, if required.

## 7.3. Assembly state detection

### 7.3.1 Implementation details

All ASD baselines make use of the YOLOv8-m [20] backbone. The baselines were trained with a learning rate of 5e-4 using the Adam optimizer, warm-up of 0.5 epoch, patience of 5 epochs, and early stopping enabled. Data augmentation is limited to HSV (with default fractions) and image scaling with a +/- gain of 0.2 for all models. For the model trained solely on synthetic data, random occlusions and image mix-up are used as additional data augmentation techniques. The random occlusions are generated using a single rectangle with random color and size, that covers at



(a) Sample one.      (b) Sample two.

Figure 6. Samples from the synthetic dataset, generated using the publicly available 3D models of all parts and Unity Perception [7].

least 50% of the bounding box and forces the class label to be background. These occlusions are generated on 33% of the training images. Image mix-up is performed by merging a synthetic training image with a random image taken from VOC2012 [12], weighing the VOC image with a random factor between 0 and 0.2.

All baselines are trained on a single Nvidia Tesla v100.

### 7.3.2 Quantitative and qualitative analysis

Two samples from the synthetic dataset constructed with Unity Perception [7], used to complement the real-world IndustReal annotations, are demonstrated in Fig. 6. Furthermore, the precision-recall curve of B3 on the test set is shown in Fig. 7. Figure 8 highlights the prediction of the B3 ASD baseline on the IndustReal test set. The figure highlights correct and incorrect predictions for challenging samples.

## 7.4. Procedure step recognition

### 7.4.1 Procedure order similarity metric

The main paper proposes to measure the quality of a predicted sequence order $\hat{y}$ (e.g., 'ACB') by comparing it to a similarity measure with respect to the ground-truth sequence $y$ (e.g., 'ABC') using string similarity. Note that this is a simplified denotation of $\hat{y}$, as it does not include the prediction time $\hat{t}_{\sigma(i)}$ and confidence $c_{\sigma(i)}$. Approaches
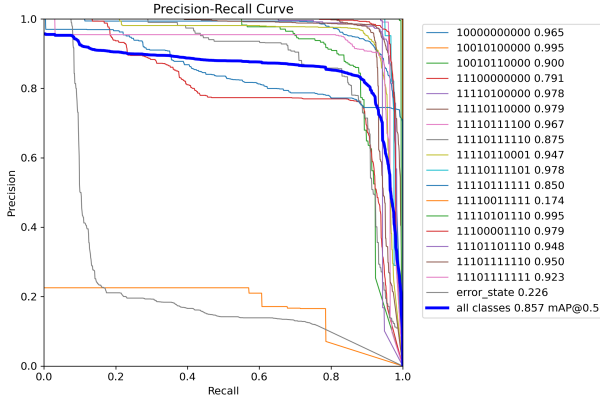
Figure 7. Precision-recall curve on ASD for YOLOv8-m [20], trained on a combination of real and synthetic data, evaluated on the IndustReal test set.

Table 6. Behaviour of the procedure order similarity (POS) metric for various predictions to the ground truth sequence $y$ "ABCD", compared to using the Lev-based similarity proposed in [23]. wDamLev indicates the weighted DamLev as proposed in Section 7.4.1, prediction errors are are underlined.

| $\hat{y}$ | AB<u>DC</u> | A<u>D</u>C<u>B</u> | <u>DB</u>C<u>A</u> | <u>_</u>BCD |
|---|---|---|---|---|
| Edits (Lev) | 2 | 2 | 2 | 1 |
| Edits (DamLev) [11] | 1 | 2 | 2 | 1 |
| Edits (wDamLev) | 1 | 3 | 4 | 1 |
| Edit (Lev) [23] | 0.50 | 0.50 | 0.50 | 0.75 |
| POS (DamLev) | 0.75 | 0.50 | 0.50 | 0.75 |
| POS (wDamLev) | 0.75 | 0.25 | 0.00 | 0.75 |

on temporal action segmentation use the Levenshtein (Lev) distance [23], normalized over the length of the ground-truth sequence . We propose to make two changes to this metric, which are outlined in more detail here than in the section in the main paper.

The edit distance is defined as the minimum number of edits required to go from a predicted to a ground-truth sequence. For the Lev distance, there are three possible edits to a sequence: insertions, deletions, and substitutions. Insertions and deletions refer to inserting or removing an element of the sequence, e.g. "ABC_" to "ABC<u>D</u>" requires the insertion of a "D". Substitutions consist of replacing an existing sequence prediction for another, e.g. "ABC<u>E</u>" to "ABC<u>D</u>" requires a single substitution of "E" to "D". The first change to the edit distance proposed in [23] is to weight substitutions with a factor of 2, rather than weighting them equally as insertions and deletions, which are weighted by a factor of 1. This essentially eliminates substitutions, as a deletion followed by an insertion is of equal edit distance. The exclusion of substitutions is proposed to prevent the metric from favouring models with many false positives. Models that falsely predict a step when it is not observed, would be penalized less than models that do not predict those false steps.

The second proposed modification is to use the Damerau-Levenshtein (DamLev) [11] edit distance, rather than the Lev distance. The DamLev edit distance allows a fourth edit method, namely transpositions, which are swaps between two subsequent elements in the sequence, *e.g.* "AB<u>DC</u>" to "ABCD". Transpositions are included in our proposed procedure order similarity, as it is intuitive to penalize "A<u>CB</u>" less compared to "<u>CA</u><u>B</u>". Transpositions, like insertions and deletions, are weighed with a factor of 1.
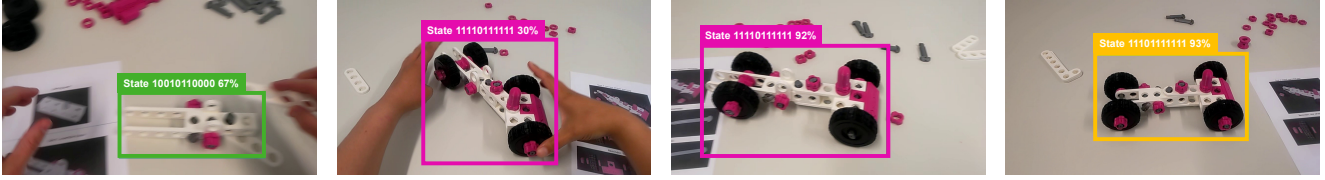
Similarly to Lea *et al.* [23], it is proposed to normalize the edit distance. The aforementioned work normalizes

with respect to the length of either the ground truth or the prediction, depending on which is longer. We propose to always normalize with respect to the length of the ground truth, as this prevents models with many false positives from being normalized favourably. Combined with the elimination of substitutions, the proposed normalization no longer bounds the normalized distance between 0 and 1. For instance, a ground truth of $(1, 2, 3)$ and a prediction of $(4, 5, 6)$ would result in a normalized edit distance of 2 (3 deletions followed by 3 insertions). Therefore, we propose to clip all normalized edit distances which are larger than 1. Since no prediction at all would yield a normalized edit distance of unity, distances larger than that are considered particularly poor predictions. Consequently, the difference between bad and worse predictions is uninteresting and not exploited in this metric. Finally, the clipped, normalized edit distance is subtracted from the unity value. This results in a similarity metric, rather than a distance metric. Thus, the procedure order similarity (POS) between $y$ and $\hat{y}$ is defined by Equation 4 in the main paper. A POS value of unity signifies a perfect match between the prediction and the ground truth ($y = \hat{y}$). The behaviour of the POS metric, as well as the DamLev edit distance, is illustrated in Tab. 6.

As mentioned in the main paper, looking exclusively at the POS metric is not sufficient for two reasons. Firstly, it does not penalize false positive predictions, if the action is later indeed correctly completed. Secondly, POS does not take time delay into account. For instance, a model could guess the order correctly before any step is completed, and obtain a perfect POS score. For these reasons, two complementary metrics are selected: $F_1$ score and the average delay $\tau$. Table 7 demonstrates the behaviour of all three metrics and outlines why the combination of them is essential.

(a) Prediction of the correct class, despite motion blur and a part occluded by the assembly.

(b) Correct prediction (with low confidence) for final state, despite a clearly misaligned front wheel assembly.

(c) Incorrect prediction (with high confidence) for the completed model, while the washer and nut for the front wheel assy are not yet completed.

(d) Incorrect prediction (with high confidence) for completed model, as the rear assy is incorrectly assembled due to a missing rear pulley.

Figure 8. Assembly state detection predictions on the IndustReal test set, using YOLOv8-m [20] (B3 in main paper).

Table 7. Demonstration of the behaviour of PSR metrics on some example predictions. A perfect POS, $F_1$ score, or delay individually does not give sufficient information regarding overall performance. The combination of the metrics provides the best insight into model quality.

| | Action \| observation time | | | | POS | $F_1$ score | Delay $\tau$ |
|---|---|---|---|---|---|---|---|
| Ground truth | $a_0 \mid 5$ s, | $a_1 \mid 10$ s, | $a_2 \mid 15$ s, | $a_3 \mid 20$ s | – | – | – |
| Prediction 1 | $a_0 \mid 5$ s, | $a_1 \mid 10$ s, | $a_2 \mid 15$ s, | $a_3 \mid 20$ s | 1.00 | 1.00 | 0.0 s |
| Prediction 2 | $a_0 \mid 5$ s, | $a_1 \mid 10$ s, | $a_3 \mid 20$ s, | $a_2 \mid 25$ s | 0.75 | 1.00 | 2.5 s |
| Prediction 3 | $a_0 \mid 5$ s, | $a_1 \mid 10$ s, | $a_3 \mid 20$ s | | 0.75 | 0.86 | 5.0 s |
| Prediction 4 | $a_3 \mid 20$ s, | $a_2 \mid 25$ s, | $a_1 \mid 30$ s, | $a_0 \mid 35$ s | 0.00 | 1.00 | 5.0 s |
| Prediction 5 | $a_0 \mid 5$ s, | $a_1 \mid 5$ s, | $a_2 \mid 10$ s, | $a_3 \mid 15$ s | 1.00 | 0.40 | 0.0 s |

### 7.4.2 Implementation details

The three baselines towards PSR that are evaluated in the paper, are further outlined in this section. Specifically, Algorithm 1 describes B1, Algorithm 2 describes B2, and Algorithm 3 describes B3. B1 consists of the most straightforward approach, where each detected assembly state that differs from the previously observed state, triggers the completion of all steps required to arrive at the detected state. B2 accumulates the prediction confidences over time, therefore filtering the ASD predictions, resulting in improved POS and $F_1$ score at the cost of an increased delay $\tau$. B3 accumulates the confidences, like B2, and additionally restricts the state transitions to the ones that are expected in the procedure. Only the expected transitions are used to predict completed procedure steps.

### 7.4.3 Qualitative analysis

Figure 9 outlines the ASD predictions on a single video in the IndustReal test set for the best scoring ASD approach, as well as the approach trained exclusively on synthetic data. These predictions are used by the PSR baselines and can therefore be used to qualitatively analyze the results from Table 4 in the main paper. It is observed that neither model is able to detect the error state around frame 1300, resulting in false positives for both PSR baselines. The model trained only on synthetic data wrongly classifies the second and last state in the video, explaining the difference in performance between B3 and B3-S outlined in Table 4 in the main paper.

**input** : List of ASD predictions $\hat{ASD}$ (in video $X_i$)
**output**: PSR predictions $\hat{y}$
$\hat{y} \leftarrow$ empty list;
$ASD_{curr} \leftarrow \hat{ASD}[0]$;
$T \leftarrow 0.5$ ;    /* Conf. threshold */
**for** $f \in \{0 \ldots \text{len}(\hat{ASD})\}$ {
    $state, conf \leftarrow getHighestPred(\hat{ASD}[f])$;
    **if** $conf \geq T$ **then**
        **if** $ASD_{curr} \neq state$ **then**
            append $state$ to $\hat{y}$;
            $ASD_{curr} \leftarrow state$;
        **end**
    **end**
}

**Algorithm 1:** PSR baseline 1

**input** : List of ASD predictions $\hat{ASD}$ (in video $X_i$)
**output:** PSR predictions $\hat{y}$
$\hat{y} \leftarrow$ empty list;
$ASD_{curr} \leftarrow \hat{ASD}[0]$;
$confs \leftarrow$ array with zeros for each component;
$T \leftarrow 8.0$ ;        /* Conf. threshold */
$decay \leftarrow 0.75$ ;    /* Confidence decay */
**for** $f \in \{0 \ldots \text{len}(\hat{ASD})\}$ {
   $state, conf \leftarrow getHighestPred(\hat{ASD}[f])$;
   **for** $i \in \{0 \ldots \text{len}(state)\}$ {
      **if** $ASD_{curr,i} \neq state_i$ **then**
         $confs_i \leftarrow confs_i + conf$;
         **if** $confs_i \geq T$ **then**
            append $state_i$ to $\hat{y}$;
            $ASD_{curr,i} \leftarrow state_i$;
         **end**
      **else**
         $confs_i = confs_i \cdot decay$;
      **end**
   }
}

**Algorithm 2:** PSR baseline 2

---

**input** : List of ASD predictions $\hat{ASD}$ (in video $X_i$), descriptive set of procedure information $\mathcal{P}$
**output:** PSR predictions $\hat{y}$
$\hat{y} \leftarrow$ empty list;
$ASD_{curr} \leftarrow \text{initialState}(\mathcal{P})$;
$confs \leftarrow$ array with zeros for each component;
$s_{exp} \leftarrow findExpectedStates(\mathcal{P})$;
$T \leftarrow 8.0$ ;        /* Conf. threshold */
$decay \leftarrow 0.75$ ;    /* Confidence decay */
**for** $f \in \{0 \ldots \text{len}(\hat{ASD})\}$ {
   $state, conf \leftarrow getHighestPred(\hat{ASD}[f])$;
   **for** $i \in \{0 \ldots \text{len}(state)\}$ {
      **if** $ASD_{curr,i} \neq state_i$ **then**
         $confs_i \leftarrow confs_i + conf$;
         **if** $confs_i \geq T$ & $ASD_{curr,i} \in s_{exp}$ **then**
            append $state_i$ to $\hat{y}$;
            $ASD_{curr,i} \leftarrow state_i$;
         **end**
      **else**
         $confs_i = confs_i \cdot decay$;
      **end**
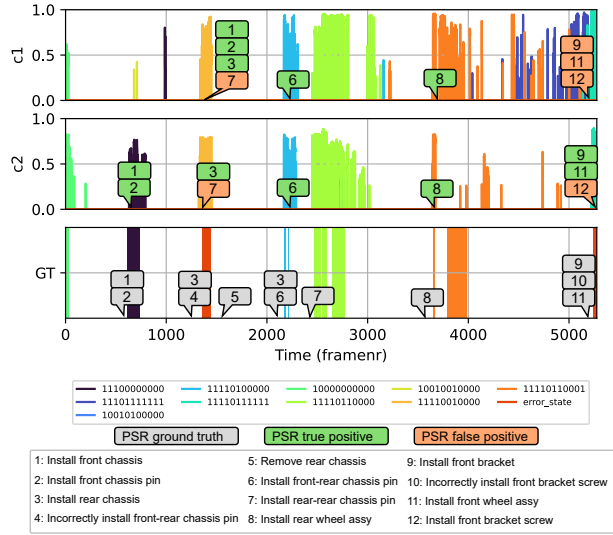   }
}

**Algorithm 3:** PSR baseline 3



Figure 9. PSR predictions by B3-S (c1) and B3 (c2), together with the accompanied ASD classifications by YOLOv8-m [20], trained exclusively on synthetic data and on a combination of real and synthetic data, respectively. Predictions are shown for a single video in the IndustReal test set. False and true positive PSR predictions are highlighted.

# References

[1] Supplementary material: IndustReal_SuppMaterial.pdf. 4, 6, 7, 8

[2] Dustin Aganian, Benedict Stephan, Markus Eisenbach, Corinna Stretz, and Horst-Michael Gross. Attach dataset: Annotated two-handed assembly actions for human action understanding. *arXiv preprint arXiv:2304.08210*, 2023. 2, 3

[3] Siddhant Bansal, Chetan Arora, and CV Jawahar. My view is the best view: Procedure learning from egocentric videos. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pages 657–675. Springer, 2022. 1

[4] Siddhant Bansal, Chetan Arora, and CV Jawahar. My view is the best view: Procedure learning from egocentric videos. *arXiv preprint arXiv:2207.10883*, 2022. 2

[5] Gregory V Bard. Spelling-error tolerant, order-independent pass-phrases via the damerau-levenshtein string-edit distance metric. *Cryptology ePrint Archive*, 2006. 4

[6] Yizhak Ben-Shabat, Xin Yu, Fatemeh Saleh, Dylan Campbell, Cristian Rodriguez-Opazo, Hongdong Li, and Stephen Gould. The ikea asm dataset: Understanding people assembling furniture through actions, objects and pose. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 847–859, 2021. 2, 3

[7] Steve Borkman, Adam Crespi, Saurav Dhakad, Sujoy Ganguly, Jonathan Hogins, You-Cyuan Jhang, Mohsen Kamalzadeh, Bowen Li, Steven Leal, Pete Parisi, et al. Unity perception: Generate synthetic data for computer vision. *arXiv preprint arXiv:2107.04259*, 2021. 7, 10

[8] Hennie Brugman, Albert Russel, and Xd Nijmegen. Annotating multi-media/multi-modal resources with elan. In *LREC*, pages 2065–2068, 2004. 9

[9] Grazia Cicirelli, Roberto Marani, Laura Romeo, Manuel García Domínguez, Jónathan Heras, Anna G Perri, and Tiziana D'Orazio. The ha4m dataset: Multimodal monitoring of an assembly task for human action recognition in manufacturing. *Scientific Data*, 9(1):745, 2022. 3

[10] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The EPIC-KITCHENS dataset. In *Proc. ECCV*, Los Alamitos, 2018. IEEE. 2

[11] Fred J Damerau. A technique for computer detection and correction of spelling errors. *Communications of the ACM*, 7(3):171–176, Mar. 1964. 4, 11

[12] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html. 10

[13] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. SlowFast networks for video recognition. In *Proc. ICCV*, pages 6202–6211, Los Alamitos, 2019. IEEE. 7, 9, 10

[14] Mingyu Fu, Wei Fang, Shan Gao, Jianhao Hong, and Yizhou Chen. Edge computing-driven scene-aware intelligent augmented reality assembly. *The International Journal of Advanced Manufacturing Technology*, 119:7369—-7381, Apr. 2022. 1, 2

[15] Carly R Garrow, Karl-Friedrich Kowalewski, Linhong Li, Martin Wagner, Mona W Schmidt, Sandy Engelhardt, Daniel A Hashimoto, Hannes G Kenngott, Sebastian Bodenstedt, Stefanie Speidel, et al. Machine learning for surgical phase recognition: a systematic review. *Annals of surgery*, 273(4):684–693, Apr. 2021. 1

[16] Tengda Han, Jue Wang, Anoop Cherian, and Stephen Gould. Human action forecasting by learning task grammars. *arXiv preprint arXiv:1709.06391*, 2017. 2

[17] Tomáš Hodan, Pavel Haluza, Štepán Obdržálek, Jiri Matas, Manolis Lourakis, and Xenophon Zabulis. T-LESS: An RGB-D dataset for 6D pose estimation of texture-less objects. In *Proc. WACV*, pages 880–888, Los Alamitos, 2017. IEEE. 1

[18] Youngkyoon Jang, Brian Sullivan, Casimir Ludwig, Iain Gilchrist, Dima Damen, and Walterio Mayol-Cuevas. EPIC-Tent: An egocentric video dataset for camping tent assembly. In *Proc. ICCVW*, pages 4461–4469, Los Alamitos, 2019. IEEE. 2

[19] Dibene J.C. HoloLens 2 Sensor Stream. *GitHub repository: https://github.com/jdibenes/hl2ss*, 2022. 6

[20] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. YOLO by Ultralytics, Jan. 2023. 7, 10, 11, 12, 13

[21] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 7, 9, 10

[22] P. Kiefe. STEMFIE construction-set toy. *https://stemfie.org/sps-000001*, 2022. 2, 5

[23] Colin Lea, Austin Reiter, René Vidal, and Gregory D Hager. Segmental spatiotemporal cnns for fine-grained action segmentation. In *European conference on computer vision*, pages 36–52. Springer, 2016. 4, 11

[24] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. MViTv2: Improved multiscale vision transformers for classification and detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4804–4814, 2022. 7, 9, 10

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 7

[26] Xudong Lin, Fabio Petroni, Gedas Bertasius, Marcus Rohrbach, Shih-Fu Chang, and Lorenzo Torresani. Learning to recognize procedural activities with distant supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13853–13863, 2022. 2, 4

14

[27] Hangfan Liu, Yongzhi Su, Jason Rambach, Alain Pagani, and Didier Stricker. TGA: Two-level group attention for assembly state detection. In *Proc. ISMAR-Adjunct*, pages 258–263, Los Alamitos, 2020. IEEE. 1, 2, 3

[28] Kosuke Moriwaki, Gaku Nakano, and Tetsuo Inoshita. The brio-ta dataset: Understanding anomalous assembly process in manufacturing. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 1991–1995. IEEE, 2022. 2, 3

[29] Francesco Ragusa, Antonino Furnari, and Giovanni Maria Farinella. Meccano: A multimodal egocentric dataset for humans behavior understanding in the industrial-like domain. *Computer Vision and Image Understanding*, page 103764, 2023. 3

[30] Francesco Ragusa, Antonino Furnari, Salvatore Livatino, and Giovanni Maria Farinella. The MECCANO dataset: Understanding human-object interactions from egocentric videos in an industrial-like domain. In *Proc. WACV*, pages 1569–1578, Los Alamitos, 2021. IEEE. 1, 2, 3, 6, 7, 9

[31] Komang Rinartha, Wayan Suryasa, and Luh Gede Surya Kartika. Comparative analysis of string similarity on dynamic query suggestions. In *Proc. EECCIS*, pages 399–404, Los Alamitos, 2018. IEEE. 4

[32] Tim J. Schoonbeek, Hans Onvlee, Pierluigi Frisco, Peter H.N. De With, and Fons Van der Sommen. Beyond action recognition: Extracting meaningful information from procedure recordings. In *2023 IEEE Conference on Virtual Reality and 3D User Interfaces Abstracts and Workshops (VRW)*, pages 881–882, 2023. 1

[33] Fadime Sener, Dibyadip Chatterjee, Daniel Shelepov, Kun He, Dipika Singhania, Robert Wang, and Angela Yao. Assembly101: A large-scale multi-view video dataset for understanding procedural activities. In *Proc. CVPR*, pages 21096–21106, Los Alamitos, 2022. IEEE. 1, 2, 3

[34] Yuhan Shen, Lu Wang, and Ehsan Elhamifar. Learning to segment actions from visual and language instructions via differentiable weak sequence alignment. In *Proc. CVPR*, pages 10156–10165, Los Alamitos, 2021. IEEE. 2

[35] Yongzhi Su, Mingxin Liu, Jason Rambach, Antonia Pehrson, Anton Berg, and Didier Stricker. IKEA object state dataset: A 6DoF object pose estimation dataset and benchmark for multi-state assembly objects. *arXiv preprint arXiv:2111.08614*, 2021. 2

[36] Yongzhi Su, Jason Rambach, Nareg Minaskan, Paul Lesur, Alain Pagani, and Didier Stricker. Deep multi-state object pose estimation for augmented reality assembly. In *Proc. ISMAR-Adjunct*, pages 222–227, Los Alamitos, 2019. IEEE. 1, 2

[37] Dorin Ungureanu, Federica Bogo, Silvano Galliani, Pooja Sama, Xin Duan, Casey Meekhof, Jan Stühmer, Thomas J Cashman, Bugra Tekin, Johannes L Schönberger, et al. Hololens 2 research mode as a tool for computer vision research. *arXiv preprint arXiv:2008.11239*, 2020. 5

[38] Zihui Xue, Yale Song, Kristen Grauman, and Lorenzo Torresani. Egocentric video task translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2310–2320, 2023. 2

[39] Xuyue Yin, Xiumin Fan, Wenmin Zhu, and Rui Liu. Synchronous AR assembly assistance and monitoring system based on ego-centric vision. *Assembly Automation*, 39(1):1–16, Apr. 2018. 2

[40] Bing Zhou and Sinem Güven. Fine-grained visual recognition in mobile augmented reality for technical support. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3514–3523, Dec. 2020. 2

[41] Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-task weakly supervised learning from instructional videos. In *Proc. CVPR*, pages 3537–3545, Los Alamitos, 2019. IEEE. 2

[42] Odysseas Zisimopoulos, Evangello Flouty, Imanol Luengo, Petros Giataganas, Jean Nehme, Andre Chow, and Danail Stoyanov. DeepPhase: surgical phase recognition in cataracts videos. In *Proc. MICCAI*, pages 265–272. Springer, 2018. 1

15