















### 3.1.1 YOLO Detection Process

The main function of this module to detect tools and pieces (e.g. nuts, bolts) and their location in the image. Although the proposed architecture is generic and any object detector could be used, in this paper the module is based on YOLO (see section 2.4.1), which provides real-time detection with high accuracy performance. A specialized dataset was designed with images of tools and parts, commonly used in the assembly processes, to carry out the learning phase. An example of the elements detected by the YOLO Detector can be seen in figure 3.

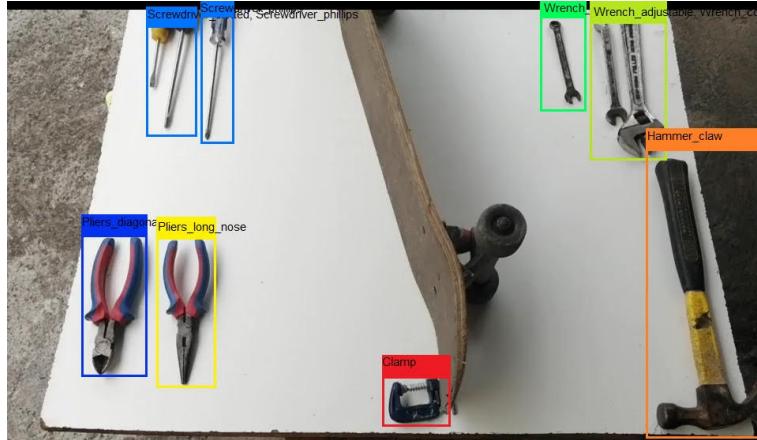


Figure 3: Example of manufacturing scene with marked objects detected by YOLO

### 3.1.2 Action Detection Process

Since *VCA*<sup>2</sup> could be mainly analyzed as an instruction recognition architecture, it must be able to identify what action the operator is performing to validate the adequate assembly or to propose the next step to carry out. This module provides this information by analyzing the sequence of video frames that come from the video capture module.

This module will generate the action as output, which combined with the object detector output will form the Action Command that will be used as input for the next stage (the Action Command Processor). The action processor engine is based on the Deep Activity Description Vector (D-ADV), a descriptor of image sequences that can describe with high accuracy the actions or activities that are being performed on it (Borja-Borja, 2020). This descriptor is the variant of the Activity Description Vector (ADV) (Azorin-Lopez et al., 2016) that can be used for deep learning methods.

The ADV and its variants have been successfully applied in different areas for the detection of actions of individuals and groups by analyzing the movement of the person as a whole. In this case, the descriptor is used to detect actions observing the working area of a production cell. The system is only able to see the tools, parts, components, and hands of the operator. The D-ADV can be used with different deep network architectures. In this paper, a two-stream architecture has been used as shown in Figure 4.

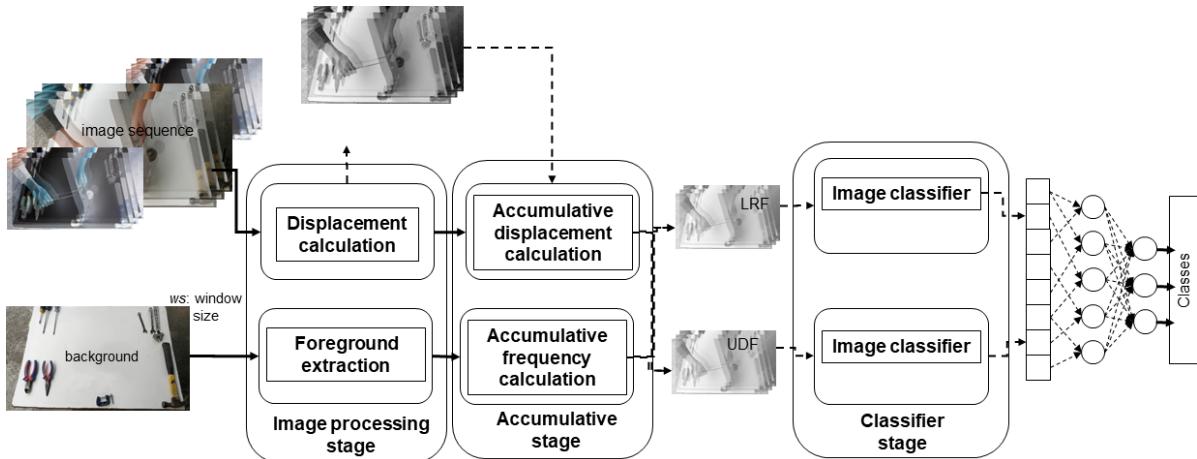


Figure 4: Deep network architecture for action detection

As seen in Figure 4, using the video that comes from the working area of the operator, the architecture **is able to calculate the action carried out by the operator in the sequence. In order to provide the corresponding class to the action, three stages are involved in the process: image processing, accumulative and classifier stages. The first two stages (image processing and accumulative stages) proceeds with the**

calculation of the image sequence representation of the deep variant of the original ADV. For the first stage, the movement occurred between two consecutive frames is calculated by means of the optical flow algorithm. The movement is separated into displacements carried out for each axis (positive or negative) considering the up (U), down (D), left (L), right (R) movements with respect to a local origin of coordinates. At the same time, the foreground of the scene is estimated to calculate the frequency (F) as the number of movements occurred in a specific region of the scene (this method divides the scenario into regions of grid cells to discretize the environment). Once the displacements and the frequency are calculated, they are accumulated according to a determined number of frames specified as the *window size* (WS). Finally, these individual movements and frequency are concatenated in D-ADV to define two images *LRF* (Left/Right/Frequency) and *UDF* (Up/Down/Frequency) which are the input queue of the classifier stage. Therefore, the D-ADV is able to compactly represent an action sequence by means of only two images. To calculate them, the input frame and the immediately WS previous frames are considered.

The classifier stage is a two-stream neural network composed of Convolutional Neural Networks (CNN) able to classify the input into many classes. The individual streams are connected at the end of the network via late fusion providing the class of the action.

### 3.2 Manufacturing description language processing

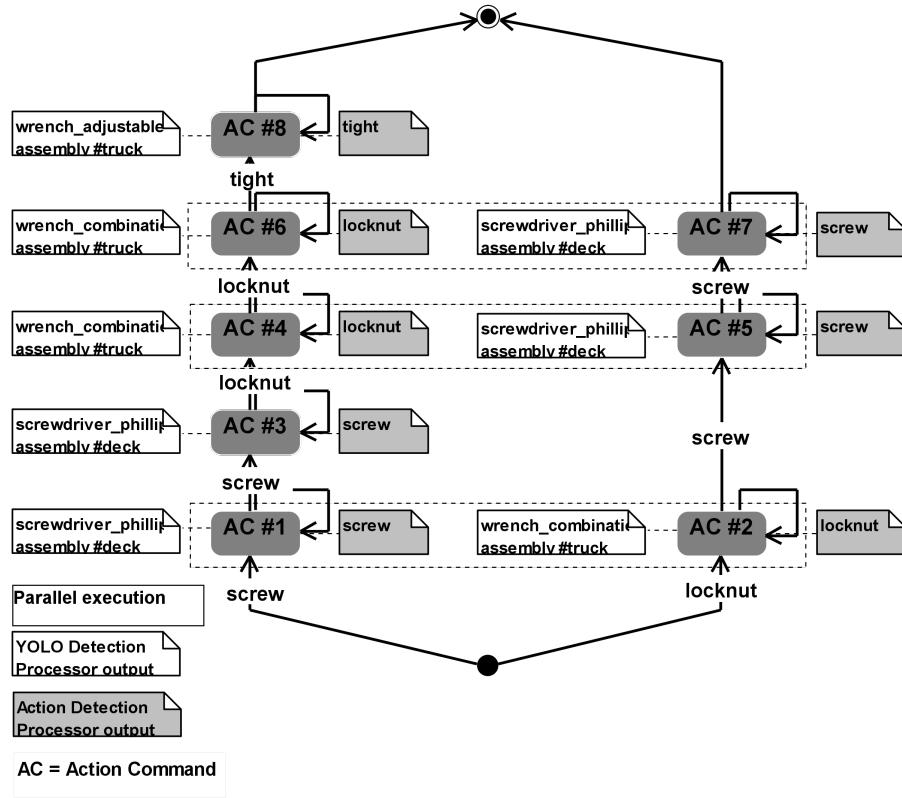


Figure 5: Example of state graph for manufacturing description language processing

This module generates a state graph with the actions, including the tools to be used, that the operator must perform (see Fig. 5. The language was designed specifically to describe operator actions (for more details see section 5.2). An example of the code can be seen in the listing 1. Consequently, it processes the assembly sequence with the instruction language, generating the expected outputs. Both the state graph and the result of the visual processing are input to evaluate the operations carried out in the working area.

Specifically, the manufacturing description language interpreter (MDLI) interprets the instructions described by the manufacturing description language to generate, as a result of the interpretation process, the state graph with the evaluated steps. Each node represents a state of the assembly and the edge, the trigger for each state change through the actions and tools to advance to the assembly process (see Fig. 6).

This graph is the one that will be used to compare the output generated from this module and that generated in the visual processing. It will allow determining if the instructions (actions and tools) were performed correctly, keeping track of the last execution carried out. With this analysis, it is possible to determine which is the current instruction to execute and the next instruction to be used as a recommendation. The interpreter generates the instructions in the format of Action Command (AC).

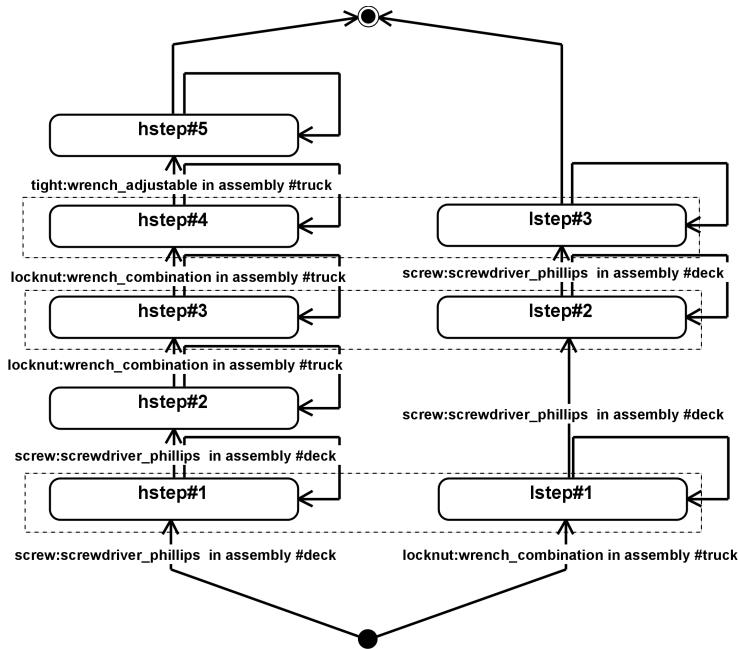


Figure 6: State diagram representing the manufacturing language that will be interpreted by the action command processor

### 3.2.1 Manufacturing Description Language for Process Control in Industry 4.0 (MDL)

The Manufacturing Description Language for Process Control in Industry 4.0 is a language that has been designed to specify the necessary actions that the operator has to perform in the assembly process. This language has the possibility of being used independently.

During the design of the language, it was established that it is a basic part of a global solution, which would use AI systems, to indicate the general operation settings automatically, such as the restrictions of the working area, tools to be used, parts and components necessary for the assembly. Anyway, for environments where this information is not available, the language provides a setup section, where the parameters can be set manually.

Another consideration of the language, during the design phase, was the syntactic structure of the instructions, which had to be simple and concise, requiring brief training in their use. Specially, since it can be used by people of different specialties, such as quality or production managers, or operators of the production cells.

Among the relevant characteristics to highlight: the method allows us to define an assembly by components, being an advantage because a milestone can be established during construction, for quality checks during assembly. It also allows establishing alternative assembly routes.

A basic set of tools was defined, but the system gives the possibility to extend with new configurations. The currently available tools are:

- |  |   |  |   |
|--|---|--|---|
| <ul style="list-style-type: none"> <li>• Clamp</li> <li>• Gun Drill</li> <li>• Screw Drill</li> <li>• Ball Pein Hammer</li> <li>• Claw Hammer</li> </ul> | <ul style="list-style-type: none"> <li>• Nut Driver</li> <li>• Diagonal Pliers</li> <li>• Lineman Pliers</li> <li>• Locking Pliers</li> <li>• Long Nose Pliers</li> </ul> | <ul style="list-style-type: none"> <li>• Ratchet</li> <li>• Electric Screwdriver</li> <li>• Phillips Screwdriver</li> <li>• Slotted Screwdriver</li> <li>• Socket</li> </ul> | <ul style="list-style-type: none"> <li>• Adjustable Wrench</li> <li>• Allen Wrench</li> <li>• Combination Wrench</li> </ul> |
|--|---|--|---|

Also as a complement to the tools, some have accessories for which they were defined in the language, as well as parts and elements of general use, among these are:

- |   |  |   |
|---|--|---|
| <ul style="list-style-type: none"> <li>• Bolt</li> <li>• Bit Drill</li> </ul> | <ul style="list-style-type: none"> <li>• Gears</li> <li>• Nut</li> </ul> | <ul style="list-style-type: none"> <li>• Screw</li> <li>• Washer</li> </ul> |
|---|--|---|



#### 4.1 ToolSet Dataset

Toolset was specifically created to satisfy our research specifications and based on the details of the proposed language called Manufacturing Description Language for Process Control in Industry 4.0 (see section 3.2.1). So, the categories of the tools are shared by both: the language and the dataset.

Toolset is an hybrid dataset consisting of 24 categories and 80,000 images including real and augmented ones. The elements contained in the dataset are general-purpose tools for manufacturing processes, tool accessories and general-purpose parts for assembly, such as screws, nuts, washers, and others.

The specifications of this dataset were designed to be trained with YOLO (see section 2.4.1), so labeling was done for this purpose, both for real pictures and for synthetic augmented images.

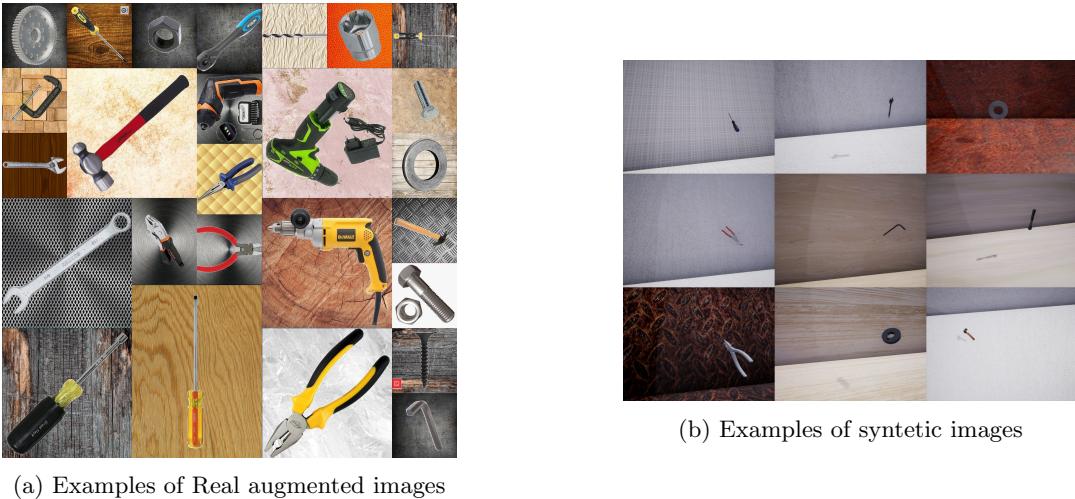


Figure 8: Examples from ToolSet Dataset

To collect the basic (real) images, we searched for each category using the datasets described in the object recognition datasets section (see section 2.3.1). But, since we did not find enough images for the defined categories, a web search was done using Google Images®<sup>11</sup> obtaining a base of 1000 real images free of use, with an intermediate to high resolution, which was later used to apply data augmentation techniques.

Regarding the data augmentation process, transformation, translation and background change techniques were applied. A total of 50 transformations were performed for each object, using 1000 images to establish random backgrounds. The generation of real data datasets takes a lot of work, since they have to search and manually tag each one of the images. To increase the size of the dataset, computational techniques were used to generate synthetic images from 3D models. To obtain a greater variety of images that would satisfy the requirements for deep network training. To generate synthetic images we used UnRealEngine, a well-known video game engine, together with a plugin called UnrealROX Martinez-Gonzalez et al. (2020), that facilitates the task of creating the synthetic dataset, allowing the creation of images from different points of view. Moreover, we performed the labeling process for the YOLO format.

#### 4.2 Video Action Dataset

To carry out the training of the action recognition architecture, we recorded several videos of manufacturing tasks to carry out the training. The camera was located on the top of the centroid of the workstation. 37 videos were recorded, from which 18 shots were obtained on average for each of the 12 actions. Figure 9 shows an example of some of the frames from one of the training videos.

The training videos were recorded in domestic and lab environments. Furthermore, recorded tools are in some cases deteriorated by use, having characteristics that usually generate noise such as wear and light oil stains. In this way, we validate detection in different environments. The videos were taken with HD web cameras, optical glass lenses, with autotracking and autofocus technologies, located one meter away from the work area.

The videos were made with simple assemblies of common elements. The assembly of a skateboard was chosen, for which different techniques and tools were used to carry out the tests. Considering different positions of the operator and different lighting levels.

To add variability to the videos, in addition to the changes made in the previous point, the sequences that describe how to build the skateboard were modified, since the assembly can be carried out in different configurations. Different positions of tools and parts within the work area were also applied.

<sup>11</sup><https://images.google.com/>







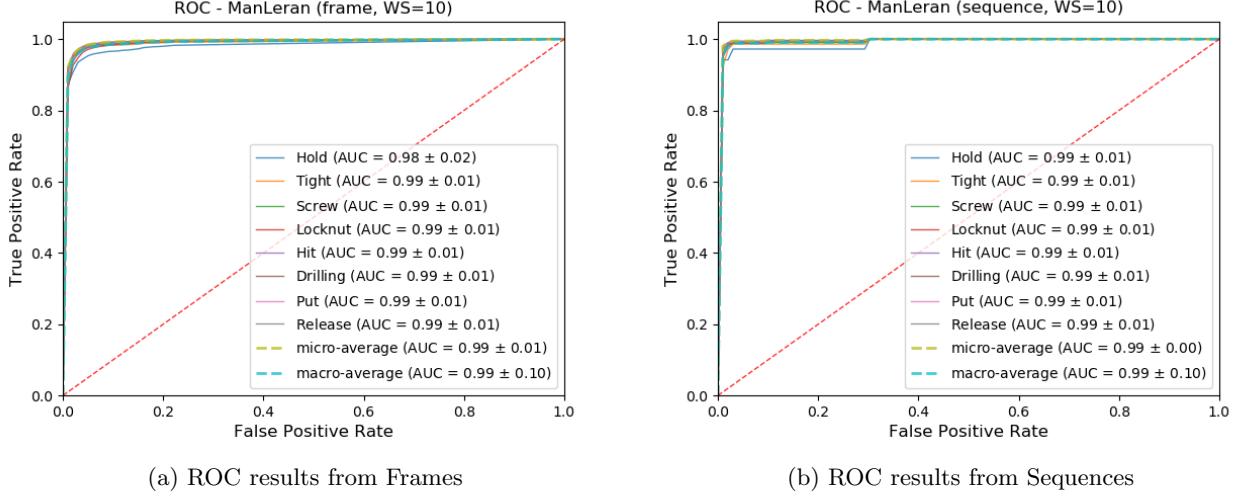


Figure 11: ROC results

This is the input together with the tokenized instructions to generate a state diagram, which is in charge of verifying that the transition between the states (a frame containing a recognized action is considered as a state). The video is reviewed in real-time to determine at which point the state changes. This state is compared with the state diagram of the instructions. If the transition is fulfilled, the process is considered to be "correct execution". The process ends when all the state transitions are completed according to the definitions until reaching the final node, the execution condition is changed to "complete execution".

All the videos tested, successfully generated the state diagrams. And therefore it was possible to generate the "complete execution" condition. The state diagram allowed us to establish which state changes were valid to make the recommendations.

In the figure 12 a representation of a process is shown, it was generated from one of the validation videos. This graph has two sections, the upper part shows the graph that is defined by the MDL instructions (the white blocks are the actions not executed and the dark ones are the actions that have already been executed). At the bottom, there is a timeline of the execution of the actions that are identified by the system, in red color is the action that is actually happening in the video and in blue color what the system identified. This is also relevant because the level of precision of detection by the system can be appreciated against reality.

The objective of the figure is to compare how the graph completes the execution of the actions in synchrony with the detection of actions by the system. In this example it can be seen that the Action Recognition Processor (ARP) in the action of the screw, detects a little earlier than the video is labeled as a screw. This is because the video is considered as a screw until the moment you turn the screwdriver, the ARP determines it from the moment the screwdriver tip touches the screw.

## 6 Conclusions

This paper presents a novel proposal to solve the prevailing problem in manufacturing for Industry 4.0 consisting of monitoring an operator during the manufacturing of a product or assembly. The proposed architecture of a visual control assistant based on Deep Learning, determines that the actions are carried out as established by the managers (quality or production engineers). Hence, the architecture is in charge of the assembly verification process and can provide systems based on it to indicate if the process was completed correctly or to alert if the steps carried out by the operator were not following the instructions. Direct benefits are obtained for stakeholders by performing this automatic assembly check: the operator improves his productivity rate, especially in those cases that are novelties or are learning the assembly of a new product, since the system can provide them with confidence; the factory improves its production capacity, reduces costs for product loss due to damage or development failures, gets shorter delivery times to customers; and, to finalize, the customer receives products in a shorter time, especially in customized mass production environments.

In this work, there were four main axes of research in order to design the architecture. Regarding the manufacturing description language, the proposal takes into account that anyone with knowledge of how a product is manufactured, be able to describe it with different levels of detail, how a product is completed.

This language provides a solution for the structured recording of the activities necessary to manufacture a product. Unlike traditional activity recording mechanisms that tend to be fuzzy and performed in natural language. Our language is structured with a general and universal syntax, without ambiguous redactions written with characteristics inherent to whoever writes them. Also, by using micro-motion theory in defining basic actions, industrial engineering experts can use it to conduct work studies. According to the experiments carried out in the paper, it was determined







- Yan, J., Yan, S., Zhao, L., Wang, Z., and Liang, Y. (2019). Research on Human-Machine Task Collaboration Based on Action Recognition. *Proceedings - 2019 IEEE International Conference on Smart Manufacturing, Industrial and Logistics Engineering, SMILE 2019*, pages 117–121.
- Yang, Y., Cai, Z., Yu, Y., Wu, T., and Lin, L. (2019). Human action recognition based on skeleton and convolutional neural network. In *2019 Photonics Electromagnetics Research Symposium - Fall (PIERS - Fall)*, pages 1109–1112.
- Yang, Y., Li, Y., Fermüller, C., and Aloimonos, Y. (2015). Robot learning manipulation action plans by "watching" unconstrained videos from the World Wide Web. *Proceedings of the National Conference on Artificial Intelligence*, 5:3686–3692.
- Zhou, L., Cao, S., Liu, J., Tan, T., Du, F., Fang, Y., and Zhang, L. (2018). Design, manufacturing and recycling in product lifecycle: New challenges and trends. *4th IEEE International Conference on Universal Village 2018, UV 2018*, 2019-Janua:1–6.