

1 **Methods of scene graph and text generation for** 1
2 **the creation of assembly instructions** 2

3 Main Track 3
4 XY 4

5 **Abstract.** In many production processes of consumer products assem- 5
6 bly instructions are often manually generated. This assembly instructions 6
7 are relevant for planning and for running a production process. In this 7
8 paper we show an approach of generating a visual graph and text outputs 8
9 with scene graph generation techniques to describe assembly instructions. 9
10 The output of a state of the art Mask-RCNN model is processed to a 10
11 Siamese Network which is combining the output information. Finally 11
12 there is applied a Graph network structure to improve the model and to 12
13 integrate context information. We compare the results to a transformer 13
14 network, which is as well producing a scene graph representation. 14

15 **Keywords:** Assembly Instructions · Scene Graph Generation · Mask- 15
16 RCNN · Graph Convolutional Network · VisionToTextTransformer · Vi- 16
17 sual Relationship Detection 17

18 **1 Introduction** 18

19 Scene graph generation is used for visual scene understanding, so that scenes in 19
20 pictures can be described with objects and their relation to each other. Objects 20
21 are described as nodes and their relation as edges. The last years there were 21
22 several approaches to solve this task. The classical approach uses a Mask-RCNN 22
23 or Faster R-CNN backbone to generate input for the actual scene generating 23
24 network. In this network semantic and visual information is processed in many 24
25 examples through MLP models to generate as output the relation of objects and 25
26 to characterize the relation. At the end of the classical approach there is often 26
27 a message passing structure integrated for example in using Graph Neural Net- 27
28 works to integrate contextual information for improving the model. The GNNs 28
29 integrate information of neighbouring nodes in a graph. 29

30 Like in many other fields transformer architectures are used to solve different 30
31 kind of machine learning problems. So as well in the context of scene graph 31
32 generation transformer technology is used to solve this task. In this approach 32
33 instead of a Mask RCNN structure there is a transformer structure like DETR 33
34 used. The output is then processed again in a further transformer architecture 34
35 to make predictions about objects and their relations to each other. 35

36 This paper shows a classical approach to solve the practical problem of gener- 36
37 ating assembly instructions with scene graph generation techniques. One objec- 37
38 tive is to show a relatively easy solution with out of the box algorithms, because 38

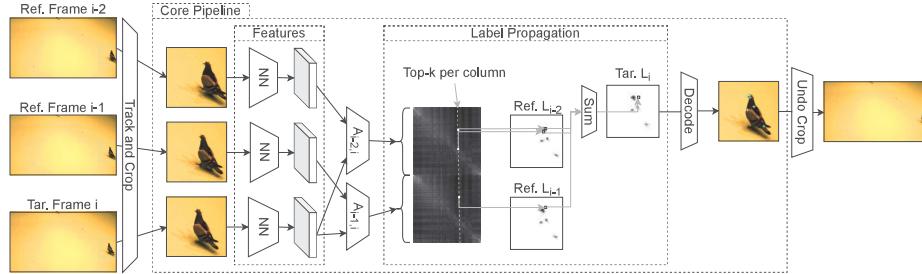


Fig. 1. Overview of the joint tracking and label propagation pipeline, showing pigeon keypoint tracking as an example. The core pipeline of our method uses deep neural network features to compute affinity matrices between a target frame and one or multiple reference frames. Labels for the target frame are computed by propagating the reference labels based on the combined affinity matrices. Finally, the predicted labels are decoded to coordinate values. The proposed extended pipeline shown outside the dashed box is necessary for tracking small objects. Here, an object tracker is used to locate the object of interest in each frame. This allows us to crop a region around the object and perform keypoint propagation at a much higher resolution, after which we undo the coordinate transform of the cropping for the decoded keypoints.

39 research paper codes are often not so easy to handle in the implementation. So
 40 as backbone network there is a Mask-RCNN architecture of MMdetection used
 41 to generate bounding boxes, segmentation masks and class prediction of objects.
 42 Then there are state of the art MLP CNNs used to model the Siamese network
 43 and finally with the library of Pytorch Geometric there is applied a GNN ar-
 44 chitecture at the end. To compare our results to the transformer approach the
 45 RelTR algorithm is applied to the same data set as the classical architecture.
 46 Another objective is to show how to solve the problem with less semantic in-
 47 formation needed. Semantic information can be an important input information
 48 for scene graph generation, because relations and the characterization of them
 49 can be derived out of semantic information. Our approach tries to put a focus
 50 stronger on geometric information of the assembly parts to make semantic in-
 51 formation less necessary. A challenge is the application of the algorithm having
 52 only a small data set, which makes generalization difficult. We use the technique
 53 of XY to make the algorithm working on a relatively small training data.

54 Main distributions of this paper are:

- 55 – To show how to implement state of the art algorithms of the classical ap-
56 proach
- 57 – To build a scene graph model relying less on semantic data
- 58 – To build a scene graph model using a Siamese network
- 59 – To use XY to handle the small dataset
- 60 – To show a specific graph structure for the practical problem

61	2 Related Work	61
62	3 Approach	62
63	The objective of the model is to be able to sufficient generalize to predict correctly an assembly instruction of unknown parts. The model is trained on toy vehicles (train, boat, plane) and should be able to predict the assembly routines of a different vehicle (car), which has similar, but unlearned parts. Building the model there is the MMdetection infrastructure needed for the Mask R-CNN network, which is trained on a dataset containig XY labeled images of the three different toy vehicles. Bounding boxes, shapes and classes of the parts are predicted. The information, if parts have a relation and if so which kind of relation is documented in graph data for training the model. The description of the graph is a crucial point. The three vehicles are combined to one graph, where different ways of modeling the graph are possible like Fig XY shows. The graph can be modeled as a vehicle based graph (Fig XY.1) or as a class based graph (Fig XY.2). We use a combination of both (Fig XY.3) to integrate vehicle specific information but as well the class information. Using an only class based graph would mean that semantic information goes more into the graph setting.	63 64 65 66 67 68 69 70 71 72 73 74 75 76 77
78	3.1 MASK R-CNN Backbone	78
79	As backbone there is used a MASK R-CNN architecture which predicts O objects in I image giving as output B bounding boxes, S shapes and C classes for each O in I. Additionally, the graph structure G is given for training. G is not represented in an adjacency matrix but as an edge list Nx2, where N is the number of nodes that have an edge connection to another node. There are $O(n^2)$ possible connections for each I, but N represents only the existing edges.	79 80 81 82 83 84
85	3.2 Siamese Network	85
86	The Siamese network is modelled in three stages in an EndToEnd structure: first the backbone output B and S is processed to the Siamese network in creating a feature vector x for B and y for S based on Resnet50 with an additional fully connected layer to allow the network to modify the extracted features. C represents the sematic information of the MASK R-CNN output and is converted to W, a word vector. X and y are then processed to individual paths, from which each is constructed in a Siamese manner. The Siamese network expects two inputs, which are processed in the same CNN with the same weights. Describing only the X path this means, that the feature vectors for two connected nodes are used by the network according the graph data G. So there is a X1 and X2 processed to the network. The Siamese architecture should allow the network to learn features comparing two images. Finally, the outputs of X1 and X2 are concatenated to one single vector X3, which is the input for a 2 layer MLP, which task is to predict the actual class of the connection between the nodes.	86 87 88 89 90 91 92 93 94 95 96 97 98 99

100 The concatenation is a direct concatenation function of the vectors from the
 101 union box of each object pair. In the example there are 4 classes used (“join”,
 102 “plug”, “screw”, “no connection”) to describe how the assembly parts can be
 103 put together. The 2 layer MLP at the end of the Siamese network combines the
 104 vectors X3, Y3 and W to predict the final class of the relation.
 100
 101
 102
 103
 104

105 **3.3 GNN Network**
 106 At the end a GNN structure is implemented to enable the learning of contextual
 107 information.
 105
 106
 107

108 **3.4 Loss function**
 109 **4 Applications and Case Studies**
 110 The generic way of formulating label propagation enables the propagation of
 111 various kinds of information through a video with a single framework. Examples
 112 include masks, keypoints and textures. Here, we take a look at mask and keypoint
 113 propagation. For each of these tasks, we have certain information in our problem
 114 domain \mathbb{P} , which we have to map to our label propagation framework. Thus,
 115 we have to define an encoding function $E : \mathbb{P} \rightarrow [0, 1]^{h \times w \times l}$ and a decoder
 116 $D : [0, 1]^{h \times w \times l} \rightarrow \mathbb{P}$.
 110
 111
 112
 113
 114
 115
 116

117 **4.1 Pose Tracking**
 118 Encoding keypoints requires some effort to transform the n points $K_{i,j} \in \mathbb{R}^2$, $j =$
 119 $1, \dots, n$ for the i th frame into 2D functions. Commonly [15,6,12,30], each key-
 120 point is encoded to a separate layer in the label function, thus $l = n$. The label en-
 121 coding for the j th keypoint is then a 2D Gaussian with mean $K_{i,j} \odot [w/W, h/H]^T$,
 122 where \odot denotes point-wise multiplication. Its standard deviation σ is a hyper-
 123 parameter and should be chosen according to the feature size. Here, we found
 124 we can improve upon some existing implementations. As the label function has
 125 the resolution of the feature map, it is often only $\frac{1}{4}$ or even $\frac{1}{8}$ of the frame reso-
 126 lution. Therefore, some inaccuracies in the placement of the Gaussian function
 127 can occur if μ is not computed with sub-pixel precision. For example, UVC [15]
 128 scales the features to the feature size but then casts the floating point numbers
 129 to integers, thus discarding the sub-pixel precision. In contrast, we implement
 130 our encoding function in such a way that the peak of the Gaussian does not
 131 necessarily lie on an integral pixel position.
 132 Decoding label functions to keypoints requires the approximate peak location
 133 for each label function. The intuitive solution of using the location of the max-
 134 imum is unreliable, as the quality of the result hinges on a single point. Hence,
 135 some methods compute the mean or center of gravity of the top k positions to
 136 improve the robustness of the solution to outliers. Our method follows the latter
 137 approach of using the center of gravity not only for its improved robustness, but
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137