# Project Proposal for ADL4CV: Scan2Cap

Felix Wimbauer
Technical University of Munich
felix.wimbauer@tum.de

Nicolas Seppich
Technical University of Munich
nicolas.seppich@tum.De

## Abstract

*Inspired by recent work on image captioning in the 2D world, we aim to investigate the task in the 3D domain. To this end, we propose a pipeline which combines concepts from 3D object detection and 2D attention-based captioning to generate a caption for a specified 3D object. The proposed pipeline first uses VoteNet to extract feature vectors of the scene. It then combines this information with the features of the object of interest, which are extracted by PointNet++, and feeds this data into an LSTM captioning mechanism, that generates a caption of the object in context of the scene.*

## 1. Introduction

The "Deep Learning" revolution has enabled great progress in the field of captioning or image description. This task has been studied in a sophisticated way for 2D images using different attention mechanisms. However, there has been no work so far to the best of our knowledge for generating a description of an object representation in point clouds or 3D data.

Therefore, we are interested in implementing a pipeline to obtain a description for a given object in a 3D scan, using state-of-the-art point cloud feature extractors, object detectors, and a captioning mechanism to generate a semantic description for a given object in the 3D scene. This allows the object to be placed in a global semantic context within its environment.

## 2. Related Work

Our work will be based on the ScanRefer dataset [2]. This dataset consists of 1513 RGB-D scans of ScanNet [3] and contains approximately 5 unique object descriptions for each object in each scene. The work of [2] will also be used as guideline in this project.

The extraction of features on point clouds is presented by [7], who apply the feature extraction directly on the point cloud on a hierarchical level, allowing the extraction of lo-cal features in a global context. The task of object detection on point clouds is studied by [6].

Methods for image captioning using visual attention are described by [9], [4] and [1]. These methods have in common, that they generate a caption for the entire image. Since our goal includes using a bounding box for the object to be set in context to the scene, the work of [8] is also of interest for this project.

## 3. Methods and Concept

Given a point cloud $p \in R^{N \times (d+C)}$ and an object in that scene, which is described by a target bounding box $b_{target} \in R^6$, our goal is to generate a meaningful caption for the object embedded in the context of the scene. To this end, we plan to use the pipeline described in the following, which was inspired by [1].

To infer information specific to the target object itself, we crop out the points belonging to the bounding box of the object and use a PointNet++ [7] network on the obtained sub point cloud. This will give a feature vector $target\_feats \in R^{128}$. In order to compute meaningful features, we use weights pretrained for point cloud classification.
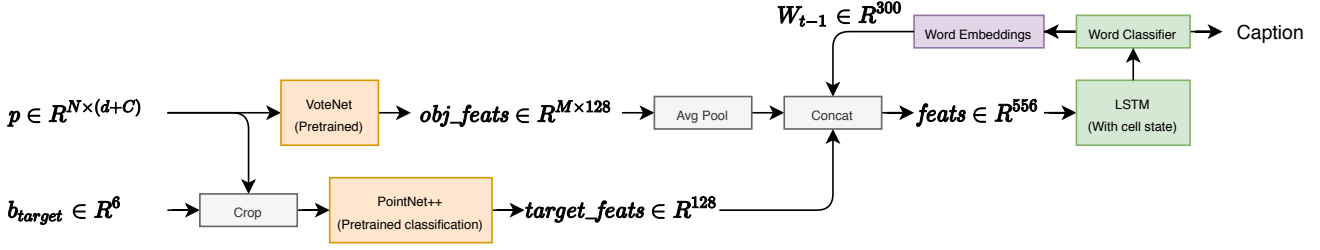
To infer information about the scene in general, we employ a VoteNet [6] network. We don't use the final object labels and bounding boxes, but the feature vectors $obj\_feats \in R^{M \times 128}$ generated by the *ProposalModule*.

In our initial pipeline, we combine those features through average pooling into $scene\_feats \in R^{128}$. At a later stage, it is possible to replace this step with attention-based pooling, similar to [1].

For the generation of the caption, we use a classical LSTM. As input, we concatenate the $target\_feats$, $scene\_feats$ and word embedding vector of the previously generated word together. The word embedding is taken from a pre-computed GloVe [5] word embedding matrix. The output of the LSTM is passed through a fully-connected layer and softmax function to obtain probabilities for the various possible next words (similar to [9]).

Figure 1 summarizes the project pipeline.

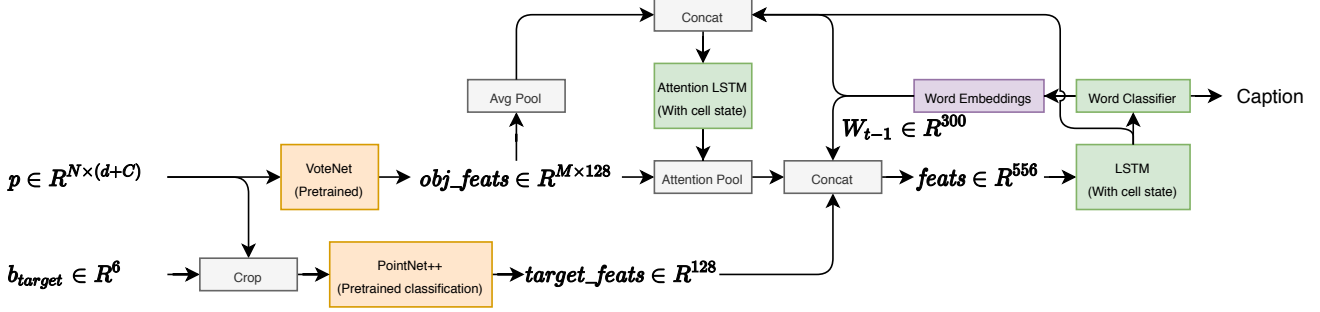**Non-Attention based captioning**



**Attention based captioning**



Figure 1. Pipeline for project

# References

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.

[2] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. *arXiv preprint arXiv:1912.08830*, 2019.

[3] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5828–5839, 2017.

[4] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 375–383, 2017.

[5] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.

[6] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9277–9286, 2019.

[7] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pages 5099–5108, 2017.

[8] Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. Grounding of textual phrases in images by reconstruction. In *European Conference on Computer Vision*, pages 817–834. Springer, 2016.

[9] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.