

# Gesture2Robot Season 2

—

## Objekterkennung über die HoloLens 2 zur Teach-in Programmierung

### **Teammitglieder**

Jan Philipp Seeland

Benedikt Beigang

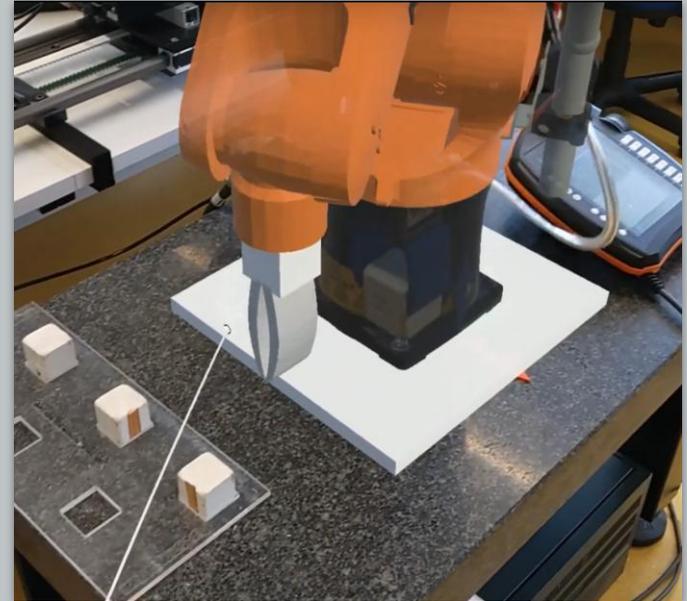
Elmar Kresse

Christoph Walther

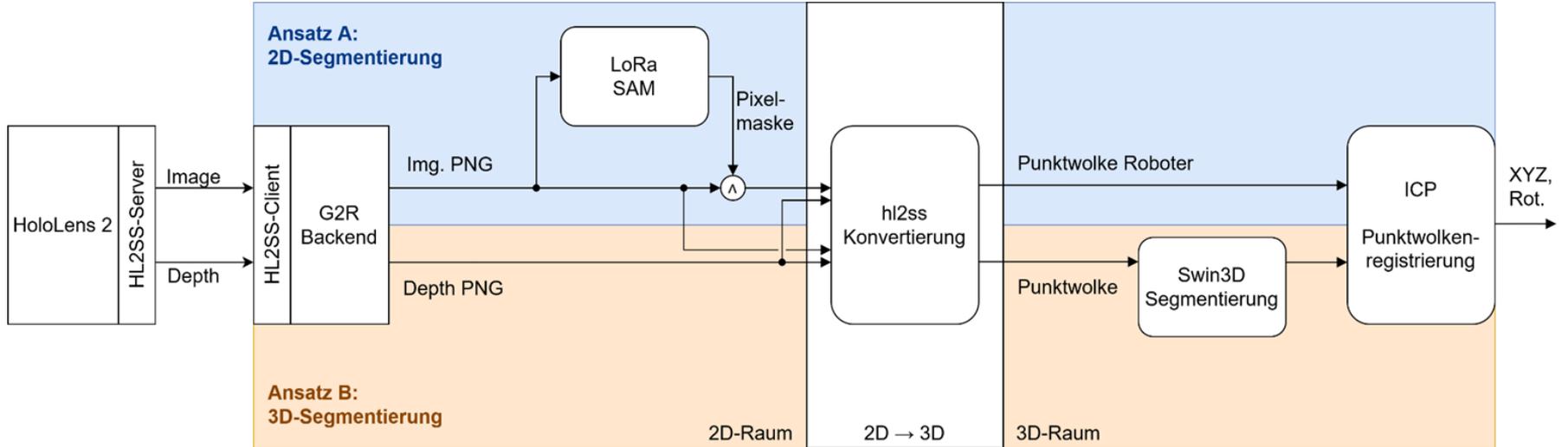
- Aufgabenstellung
- 
- 2D Segmentierung
- 
- 3D Punktwolken Segmentierung
- 
- Fazit und Zukunft

# Aufgabenstellung

## Automatisierte Platzierung des Roboters



# Zwei parallele Lösungsansätze



# Datenabfrage, Aufarbeitung und Verarbeitung

- visible Light Cameras - Four grayscale cameras 640x480@30FPS Stream
- RM Depth Long-throw - 32bit PNG 320x288@5FPS
- RGB camera (PV) - 1920x1080@30FPS Stream
- Spatial input - Head pose, gaze ray (origin and direction), hand tracking



# Datenabfrage, Aufarbeitung und Verarbeitung

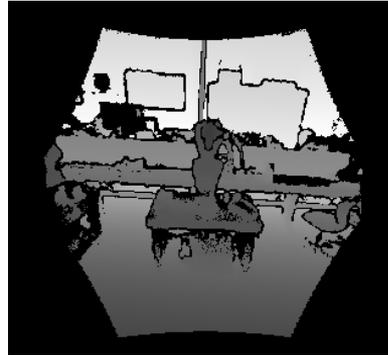
- Forschungsprojekt am Stevens Institute of Technology, gesponsert von der Defense Advanced Research Projects Agency (DARPA)
- Auswahlgrund: Bidirektionale Kommunikation + Unity-Integration auf Brille
- Motivation: Bereitstellung von Hololens-Rohdaten für rechenintensive Applikationen
- Kein Zugriff auf Tiefendaten über Microsoft-Schnittstellen  
hier: Research Mode

Field	Type	Length in bytes
Focal length	1x2 float	8
Principal point	1x2 float	8
Radial distortion	1x3 float	12
Tangential distortion	1x2 float	8
Projection	4x4 float	64

Table 15. PV calibration data.

# Datenabfrage, Aufarbeitung und Verarbeitung

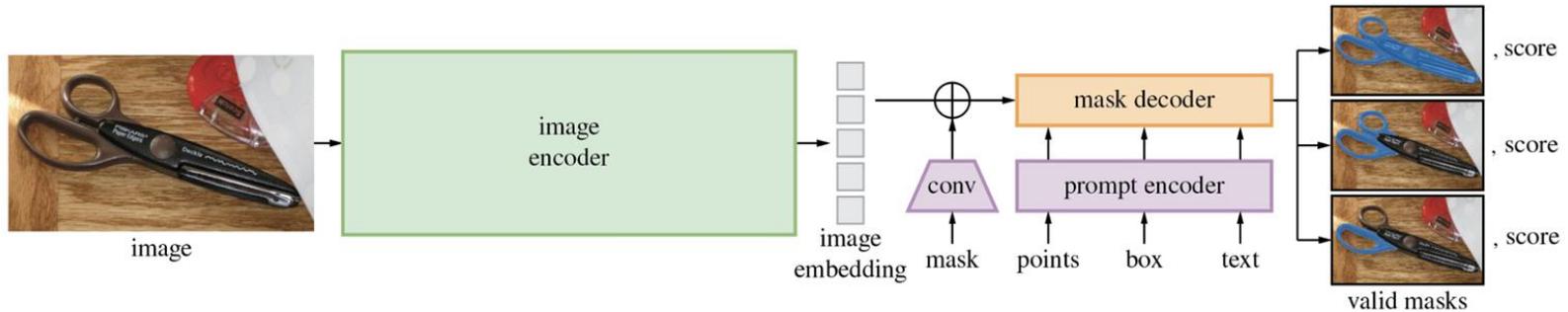
- RGBD Bildausrichtung der Aufnahmen  
Zuschnitte aus beiden Sensoraufnahmen
- Farbe kann mithilfe des research Modus und der Kameras für sichtbares Licht (oben) erzeugt werden
- alternativ Erzeugung mithilfe der Foto/Video-Streams unten
- Video + extrinsische Kalibrierungsdaten
- Autofokus deaktiviert, damit Daten valide bleiben



# 2D-Segmentierungs Ansatz **SAM** mit LoRa



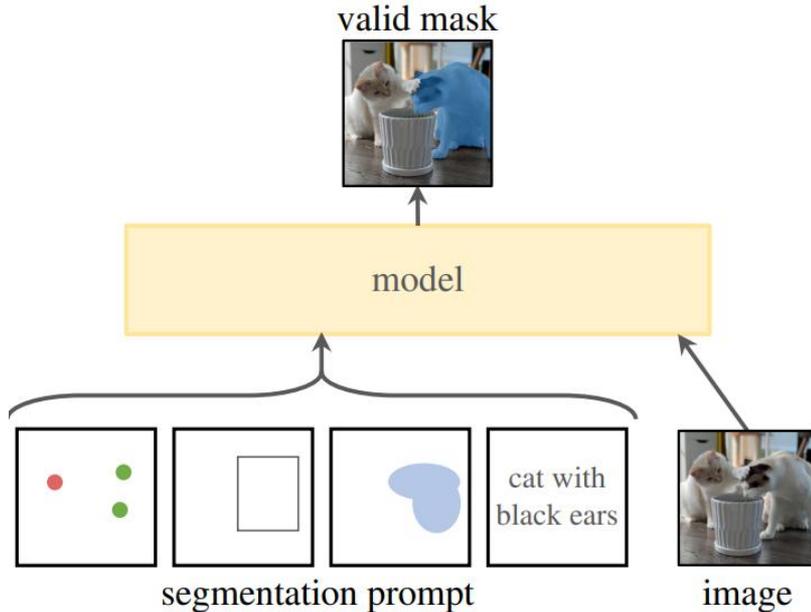
# SAM - Segment Anything Model



- trained on
  - a dataset of 11 million images and
  - 1.1 billion high quality object masks
- strong zero-shot performance

[a]

# SAM - Segment Anything Model



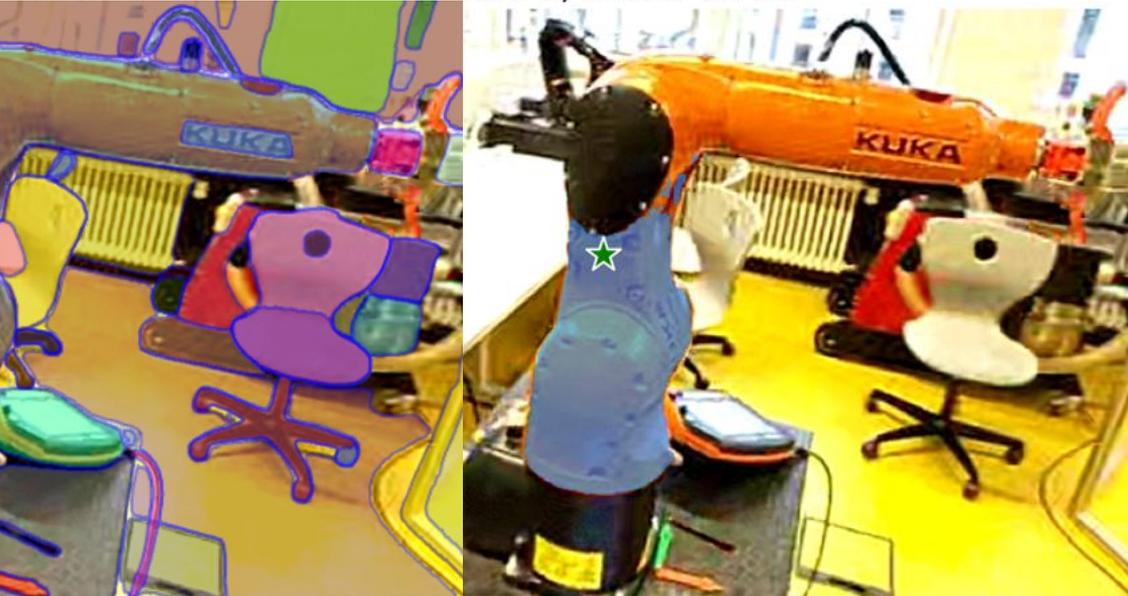
## Promptable Segmentation

- Punkte
- Bounding-Boxen
- Masken
- Text Prompt (noch nicht offiziell veröffentlicht) [9]

[a]

# SAM - Sieht doch gut aus?

Mask 0, Score: 0.719



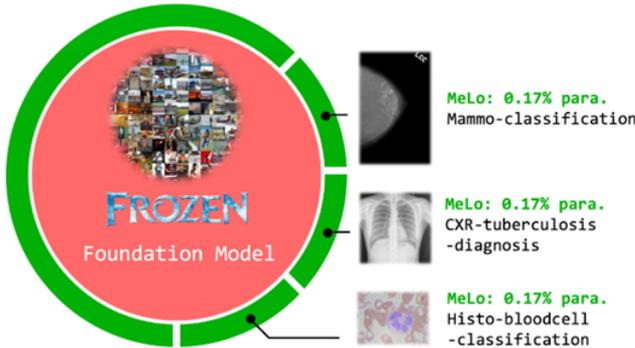
## Probleme:

- Erkennung des Konzeptes Roboter
- Geschwindigkeit: ~3-5s pro Bild
- Größe (2,4GB!)
- Finetuning nicht offiziell unterstützt

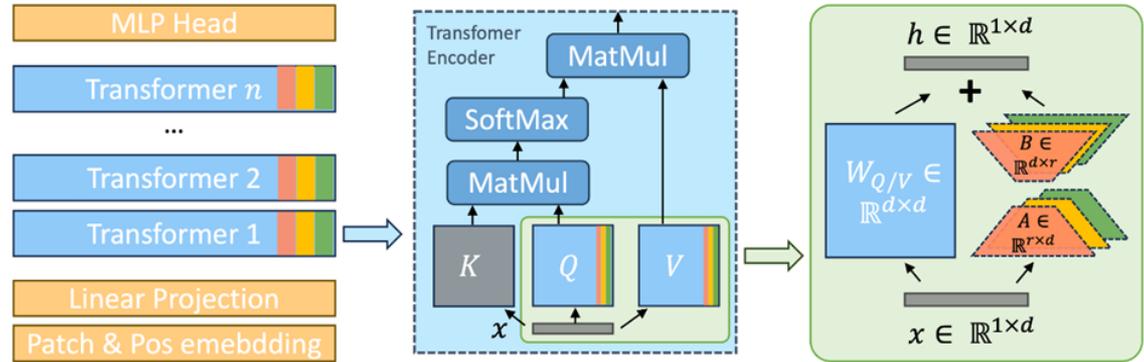
# Finetuning Ansätze

- **Segment Anything Model with 🤗 Transformers** [h]
  - sehr aktuell
  - Finetuning für Zellen und Organe
  - Problem: Training verschlechtert Ergebnisse (vermutlich zu kleiner Datensatz)
- **Segment Anything in KerasCV** [i]
  - Maskenerstellung durch Text-Prompt
  - kompliziertes SetUp (jax-backend + GroundingDINO + SAM)
  - nicht für Finetuning ausgelegt
- **LoRA-SAM (Low-rank Adaptation)** [j]
  - Ausgelegt auf fine-tuning
  - schnell im Training und beim Abspielen
  - kleine Dateien

# LoRA - Low-Rank Adaptation



(a) Overview



(b) Method

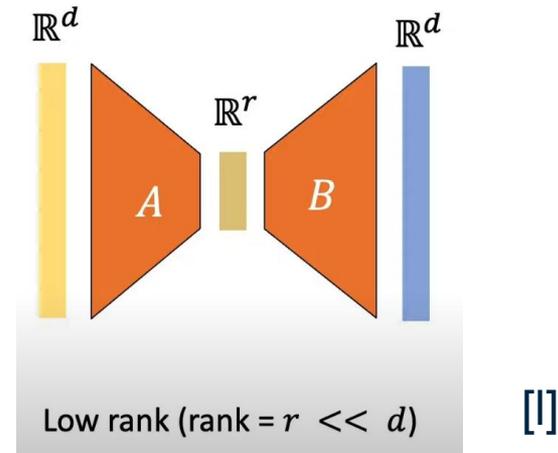
[k]

- Parameter effiziente fine-tuning Technik
  - nur wenig Anpassung der Gewichte auf vor-trainierten Modellen nötig
- Eingabe sind “low-rank decomposition” Matrizen
- erstmals in LLMs verwendet -> Übertragung auf Vision Transformer

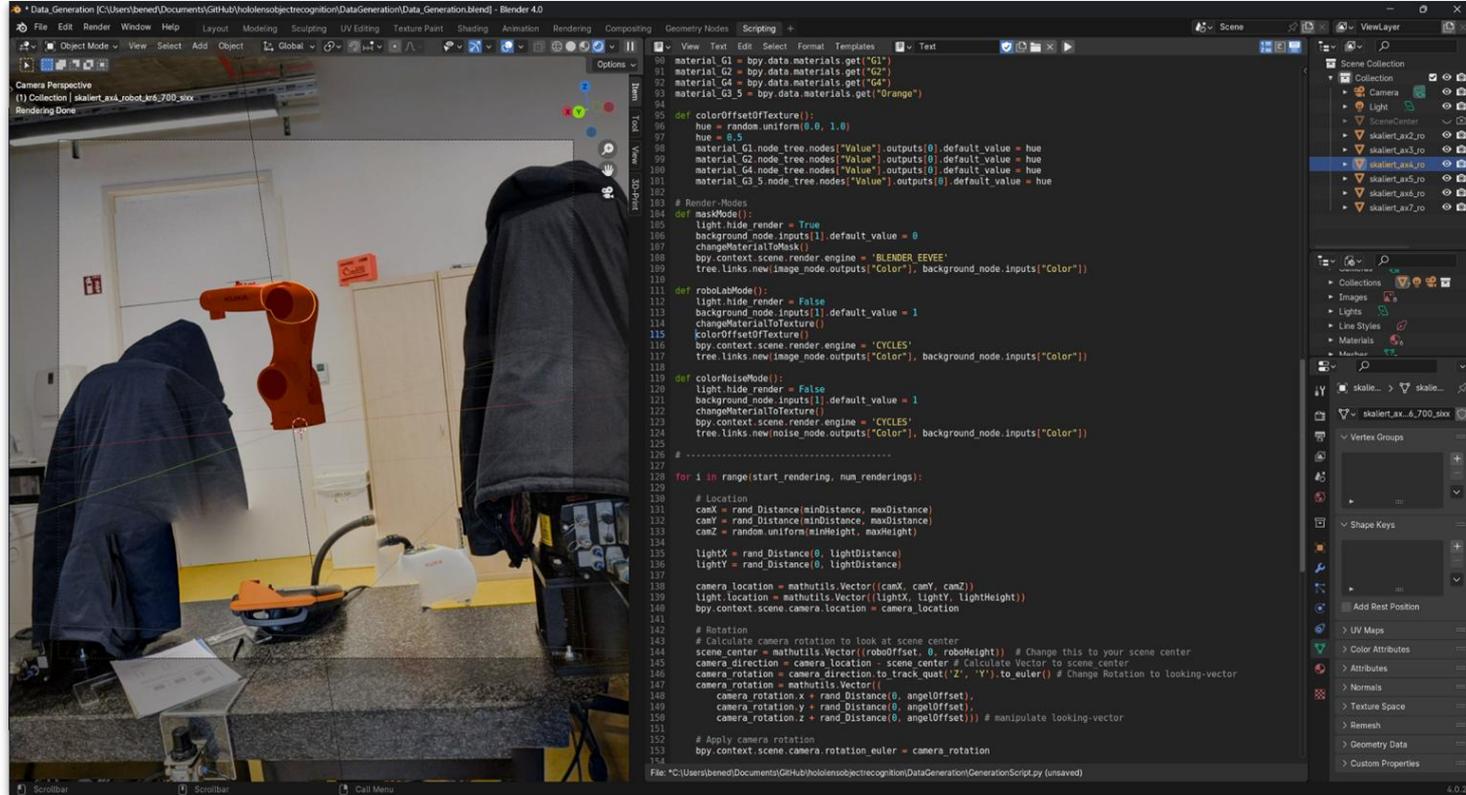
# Vorteile von LoRA

- SAM => 4,4s
- SAM mit LoRA => 0.7s
- Kleinere Datensätze
- Finetuning über kleine Ränge

*“[LoRA] freezes the original weights of visual foundation models while adding small low-rank plug-ins that can achieve remarkable results using only a small fraction of trainable weights.”*



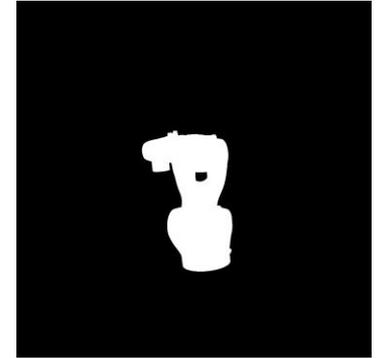
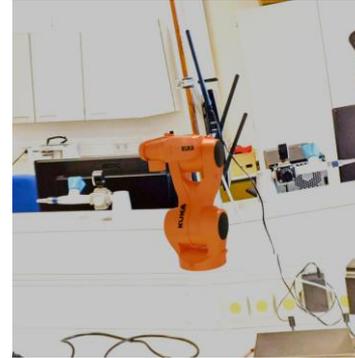
# Erstellung des Datensatzes in Blender



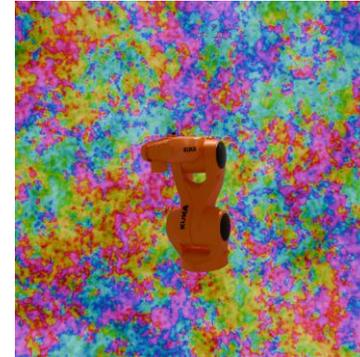
Benedikt Beigang, Jan Philipp Seeland, Elmar Kresse, Christoph Walther

# Variable Parameter

- Kameraposition
- Kamerawinkel
- Farbe des Roboters
- Beleuchtungswinkel
- Hintergrund

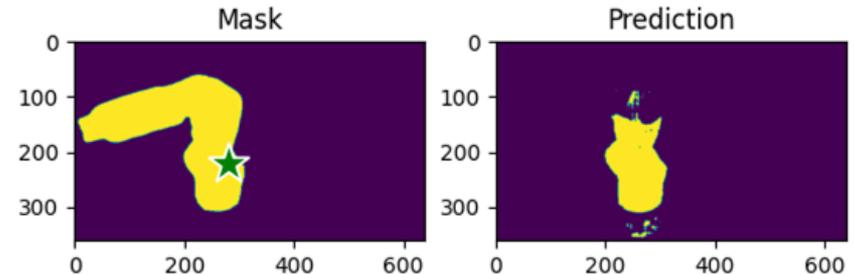


+



# Erstellung des Evaluations-Skripts

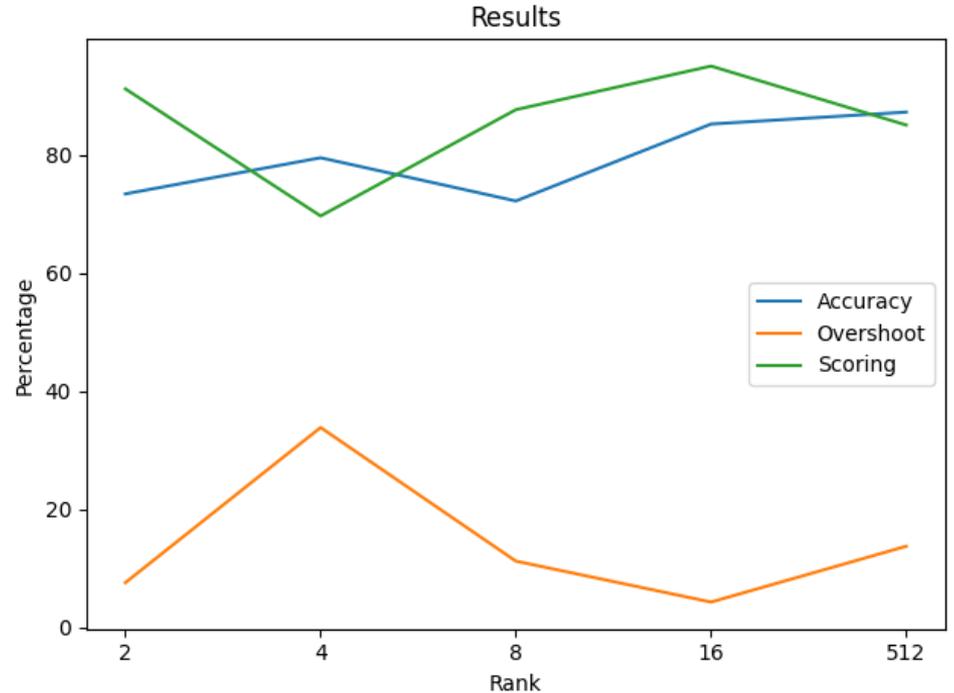
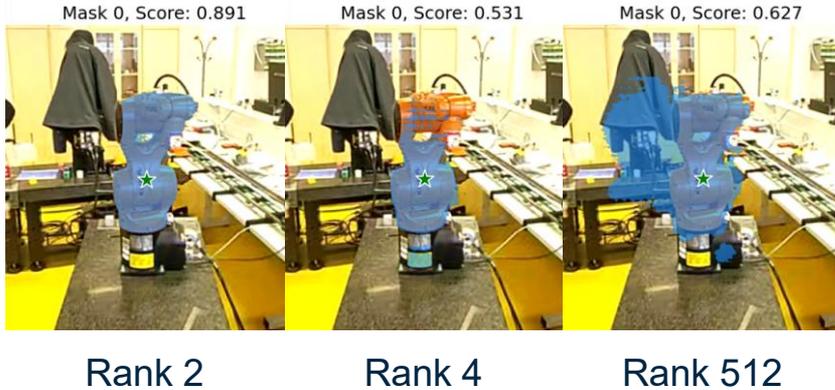
- 30 diverse Bildaufnahmen von der HoloLens 2
- Erstellung der Groundtruth-Masken per Hand
- Berechnung der falsch-negativen und falsch-positiven Pixel zur Berechnung der Accuracy und des Overshoots
- Score:  $1 - ((fp+fn) / \text{alle Pixel})$



Mask-Accuracy: 40.07 %  
Background-Overshoot: 0.32 %  
Final Scoring: 70.0 %

# Welchen Rank?

- 500 "realistisch" gerenderte Images
- 5 Epochen
- => Rank 2 und 16 am besten



# Welcher Datensatz?

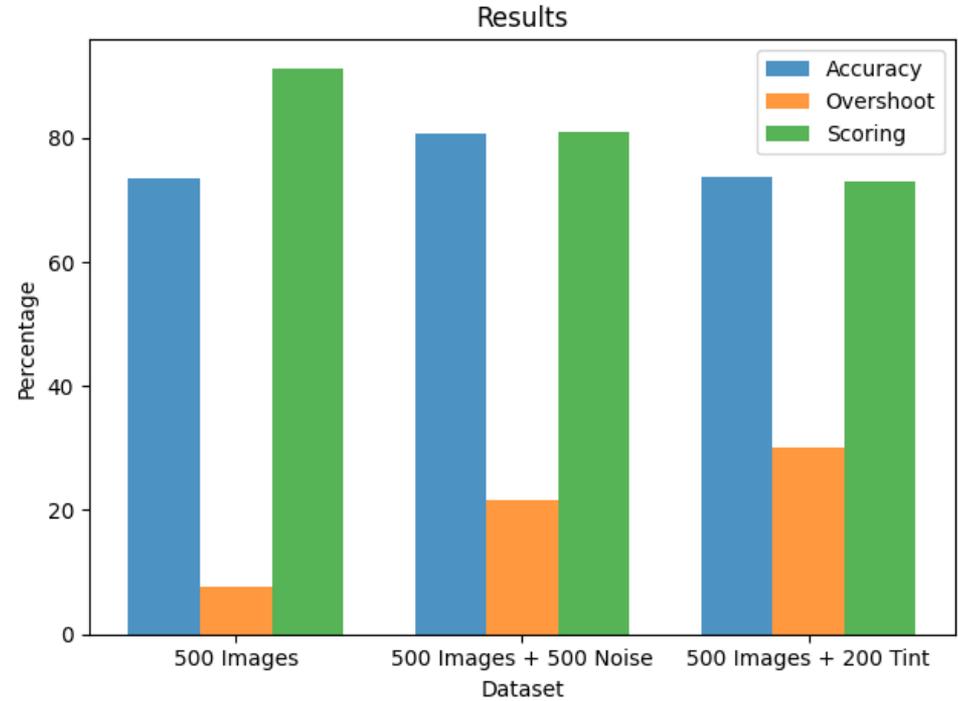
- Rank 2
- 5 Epochen
- “realistisch” gerenderte Images performen am besten



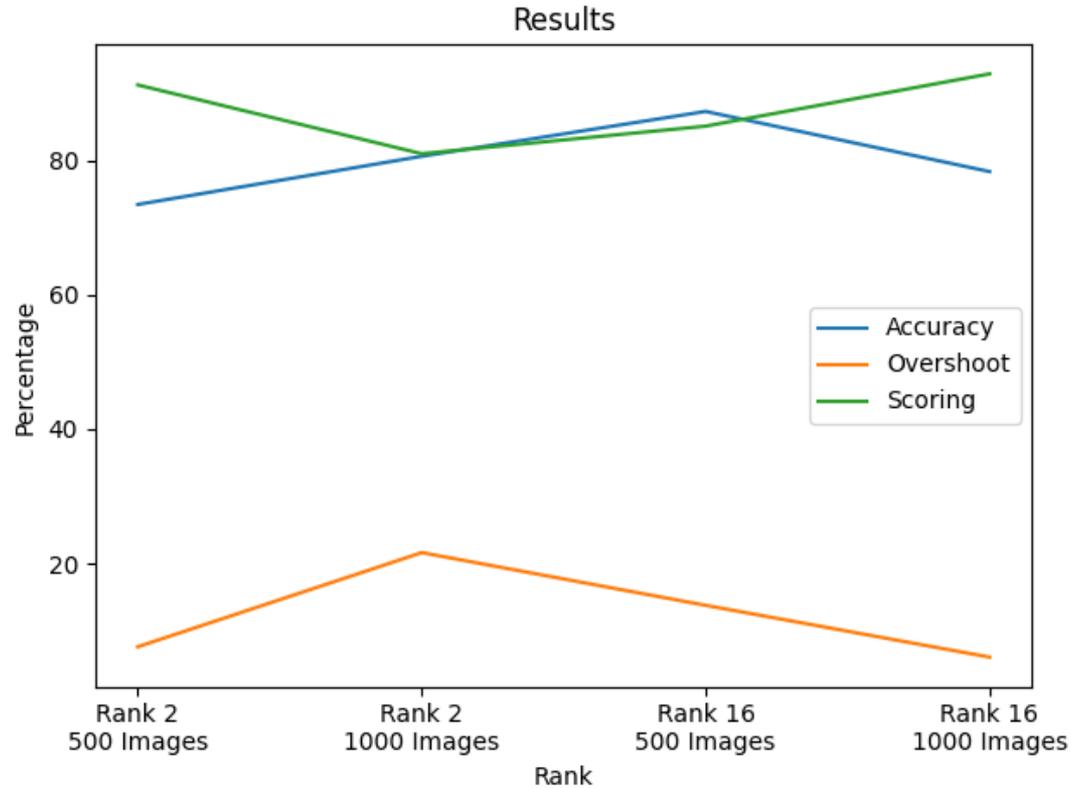
Image

Noise

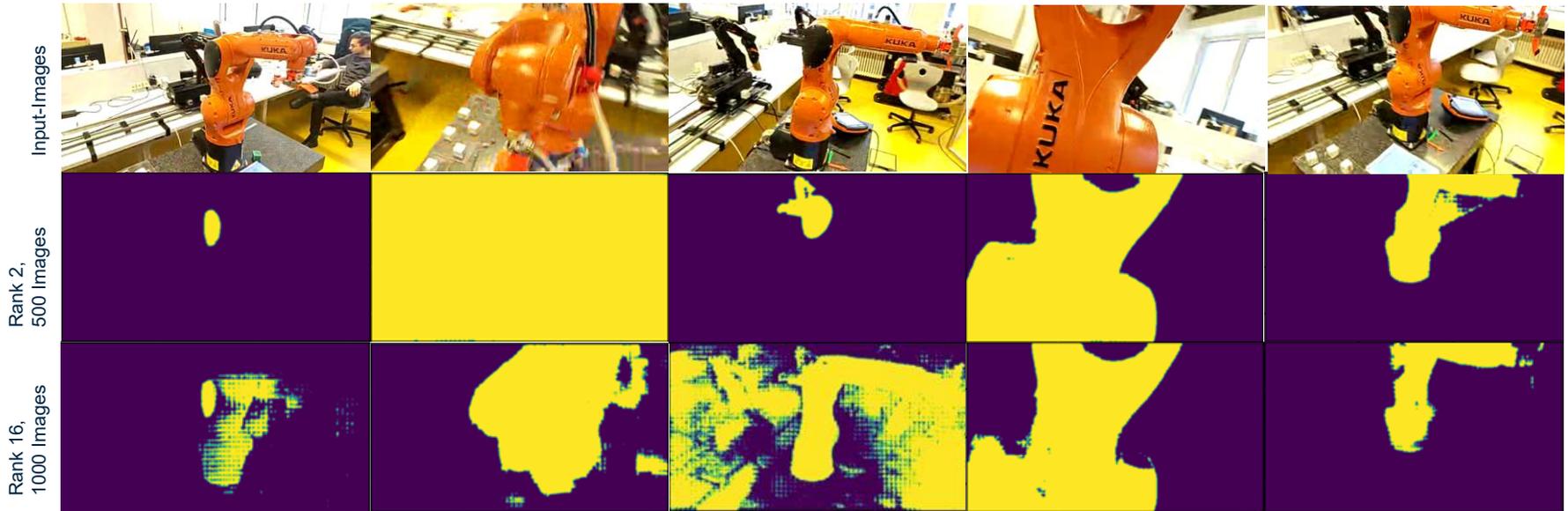
Tint



# Die besten Ergebnisse

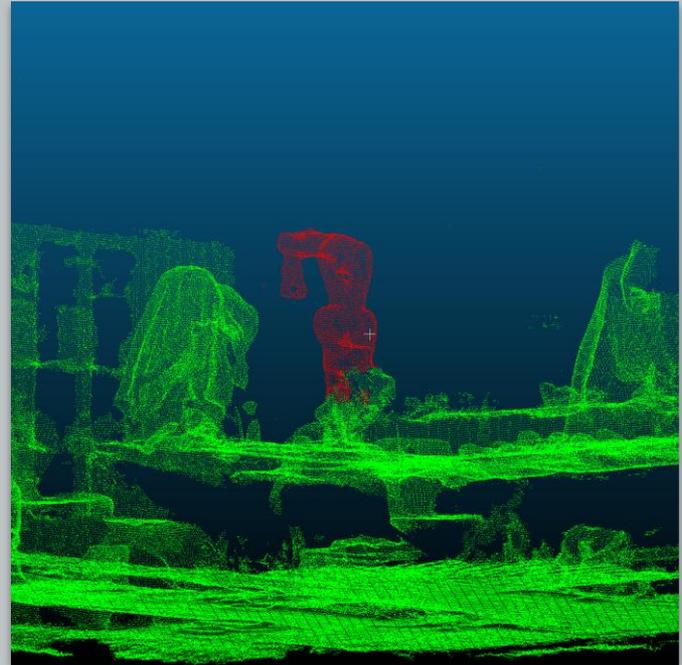


# Erfolge und Probleme

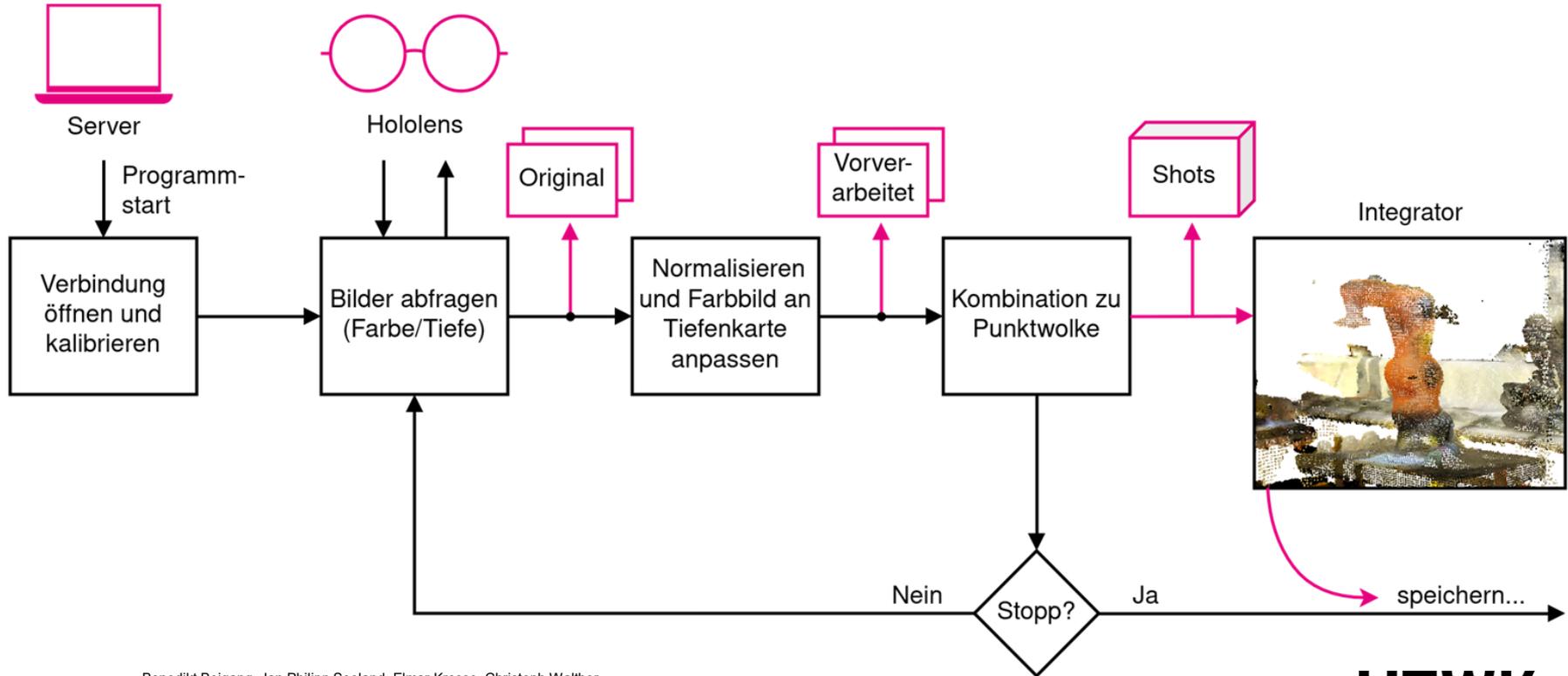


# 3D RGB Punktwolken- Segmentierungsansatz

## **Swin3D** **shift-window based** **transformer**



# PV Integration Pipeline



# Stanford 3D Indoor Scene Dataset (S3DIS)

- 6 Innenraumbereiche (“Areas”)
- Rund 70.000 Bilder, 700 Mio. Punkte
- 13 Objektklassen

## Aufbau:

- Area\_1
  - conferenceRoom\_1
    - Annotations:  
beam, board, chair, ...
  - conferenceRoom\_2
  - ...
- Area\_2...6



# Generierung eigener Datensatz

## HTWK Robo 3D Dataset (HR3D)

- Shots\_00
  - scan\_00...24: **Roboterscans**
    - Annotations: robo, other
  - scan\_50...51: **Raumscans**
- Shots\_01...24
  - pcd\_robo\_xx\_01...343

**Einzelperspektiven**
- Shots\_50...51
  - pcd\_robo\_5x\_01...180

**Einzelperspektiven Raum**
- Shots\_666: S3DIS Area 1



# Labelling der Daten (Ausschneiden des Roboters)

## 1. Manuelles Labeling

- CloudCompare:  
2D Schnitt mehrerer Blickwinkel
- Für alle 26 integrierten Punktwolken



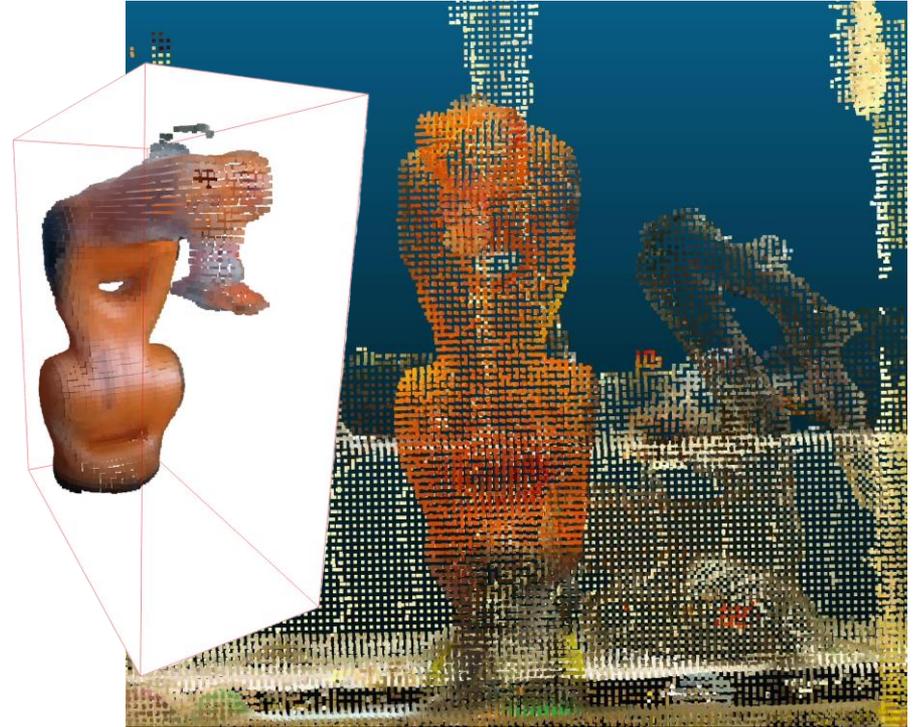
# Labelling der Daten (Ausschneiden des Roboters)

## 1. Manuelles Labeling

- CloudCompare:  
2D Schnitt mehrerer Blickwinkel
- Für alle 26 integrierten Punktwolken

## 2. Automatisches Labeling

- Bounding-Box um Roboter erstellen
- Alle Teilaufnahmen im Bereich segmentieren



# Datenaugmentierung

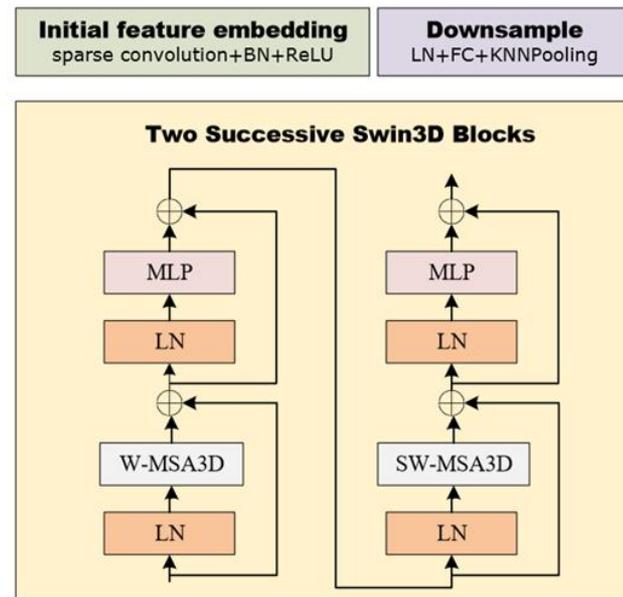
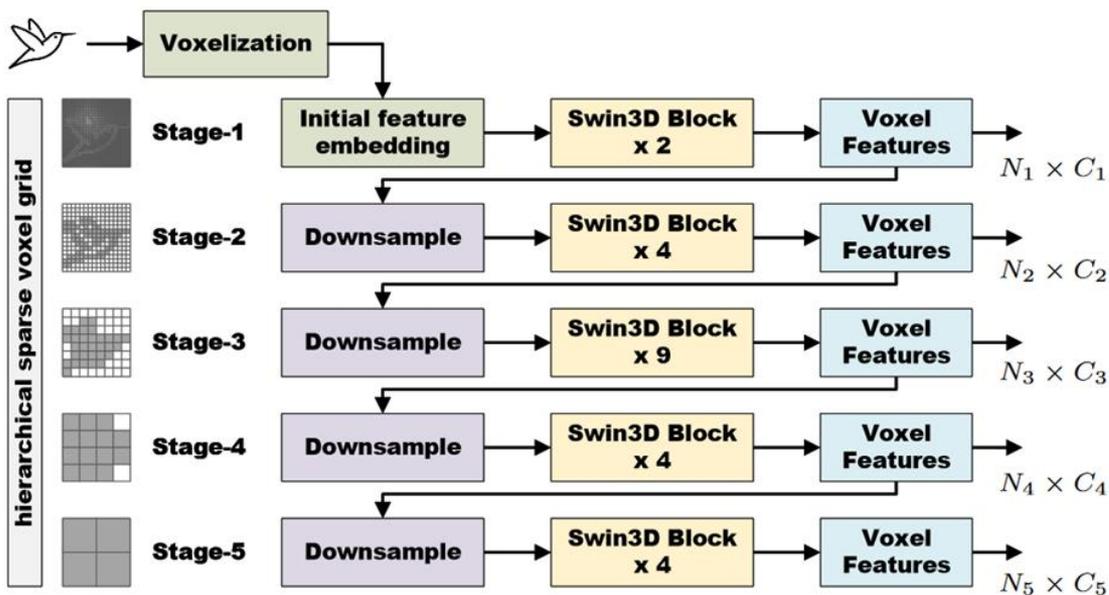
- Variierung der Scans
  - Verschiedene Roboterpositionen, Lichtsituationen, Scanqualität

## In Swin3D

- Etablierte Datenaugmentation as Papern CAGroup3D und FCAF3D
- Normalisierung und Anwendung verschiedener Transformationen
- Weitere Transformationen vorhanden (z. B. Farbe, etc) → Zukunft

```
train_transform = transform.Compose(  
    [ ...  
        transform.RandomRotate(along_z=args.get("rotate_along_z", True)),  
        transform.RandomScale(  
            scale_low=args.get("scale_low", 0.8),  
            scale_high=args.get("scale_high", 1.2),  
        ),  
        transform.RandomJitter(sigma=jitter_sigma, clip=jitter_clip),  
        transform.RandomDropColor(  
            color_augment=args.get("color_augment", 0.0)  
        ),  
    ]  
)
```

# Swin3D

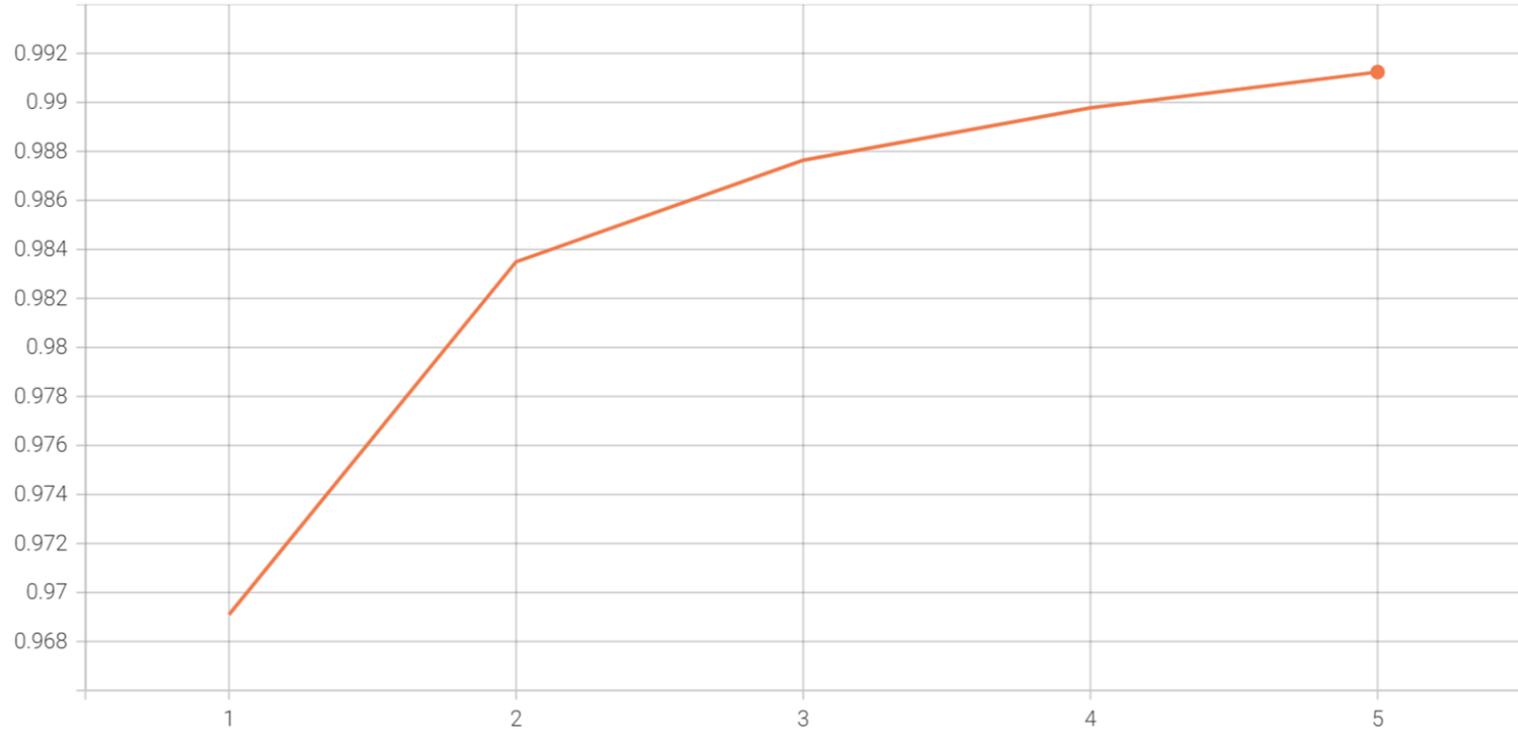


[e] Swin3D: A Pretrained Transformer Backbone for 3D Indoor Scene Understanding



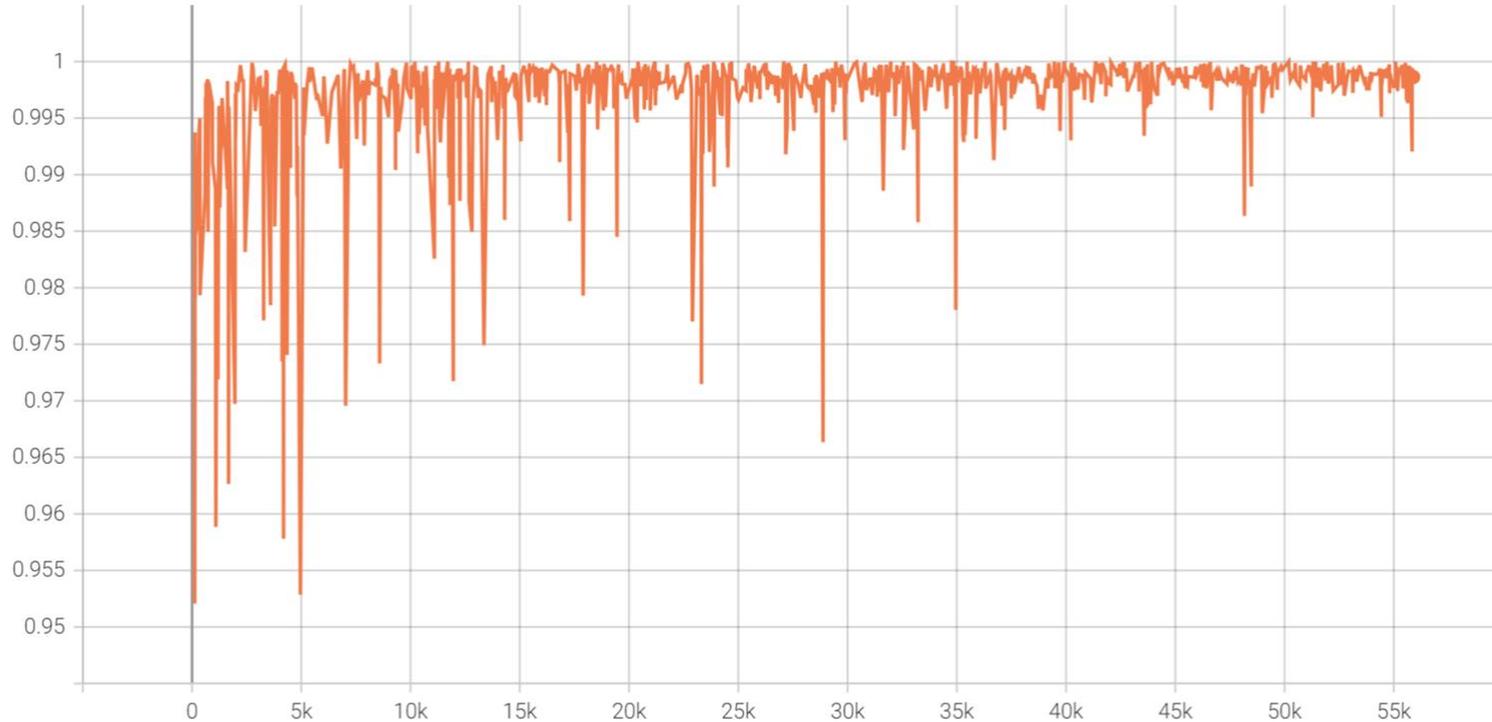
# Auswertung Training und Evaluation

mAcc\_train  
tag: mAcc\_train



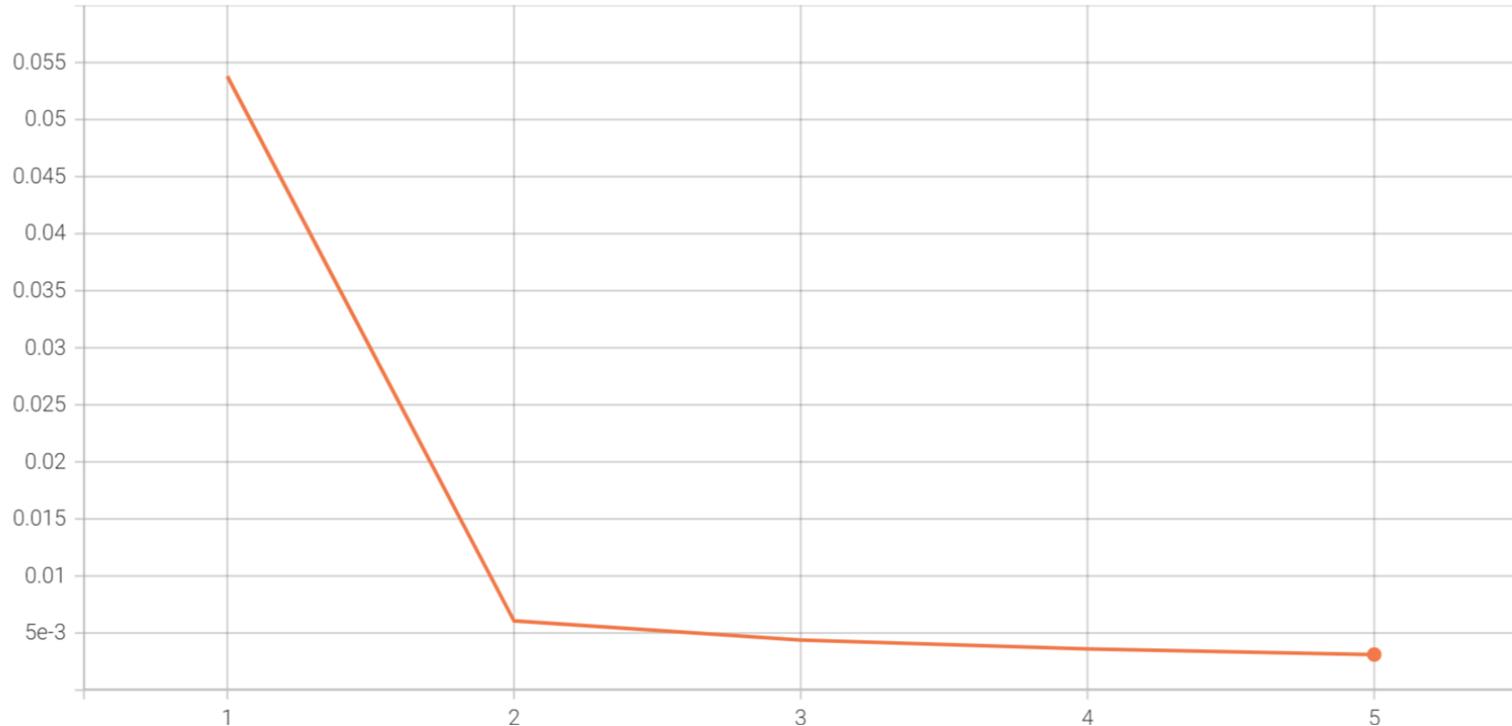
# Auswertung Training und Evaluation

allAcc\_train\_batch  
tag: allAcc\_train\_batch



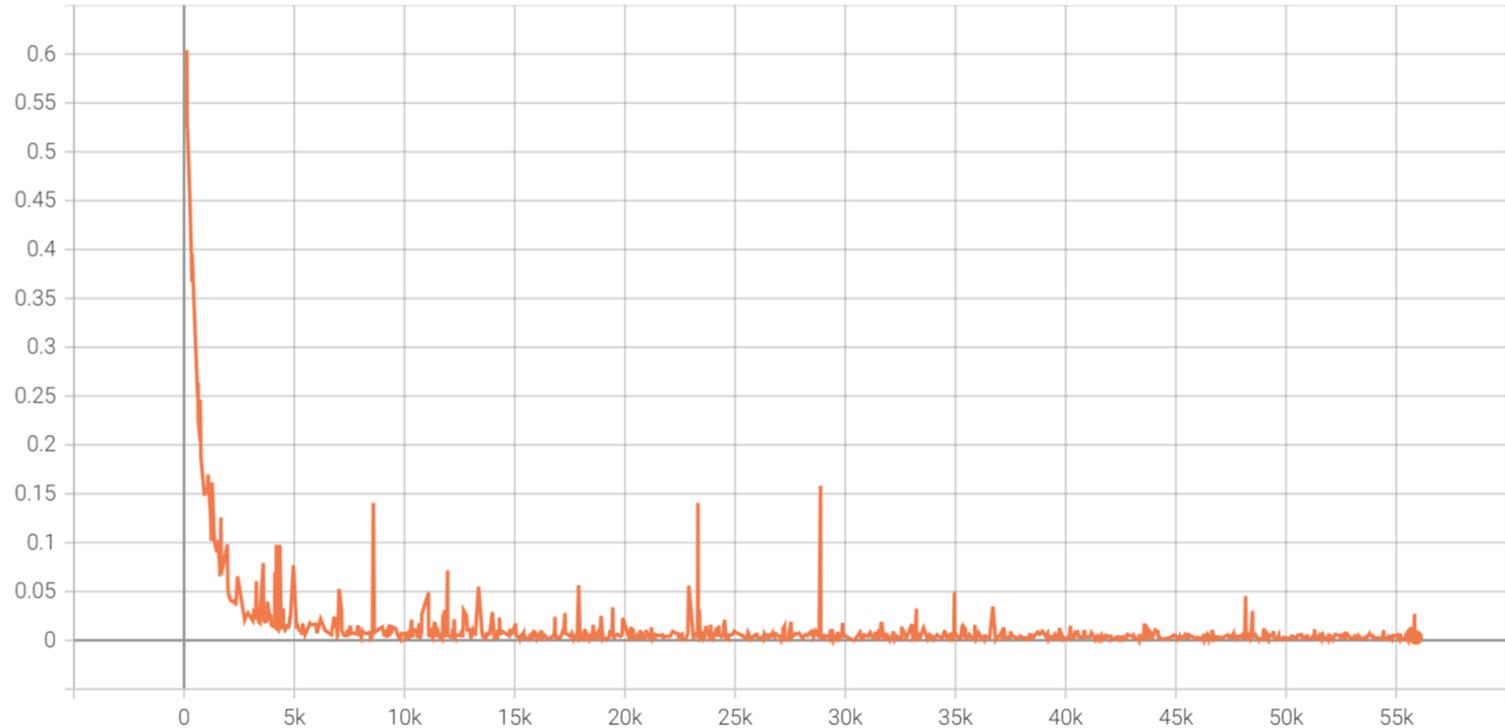
# Auswertung Training und Evaluation

loss\_train  
tag: loss\_train



# Auswertung Training und Evaluation

loss\_train\_batch  
tag: loss\_train\_batch



## Point Cloud Registration:

- Definition: Der Prozess der Ausrichtung mehrerer Punktwolken, um eine konsistente Referenz zu schaffen
- Anwendungen: 3D-Rekonstruktion, Objekterkennung, Umgebungs-Modellierung
- Herausforderungen: Rauschen, Ausreißer, nicht-überlappende Bereiche, große Datenmengen

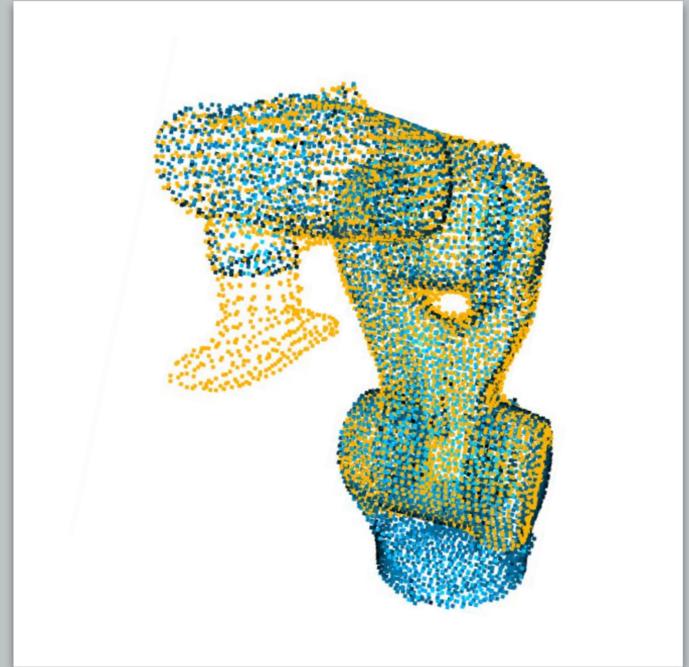


## ICP (Iterative Closest Point):

- Grundprinzip: iterativer Algorithmus zur Bestimmung der bestmöglichen Transformation zwischen zwei Punktwolken
- Ziel: Minimierung quadratischer Abstände zwischen den Punkten der beiden Wolken
- Anwendungen: Registrierung von 3D-Modellen in Robotik, Computer Vision, Augmented Reality

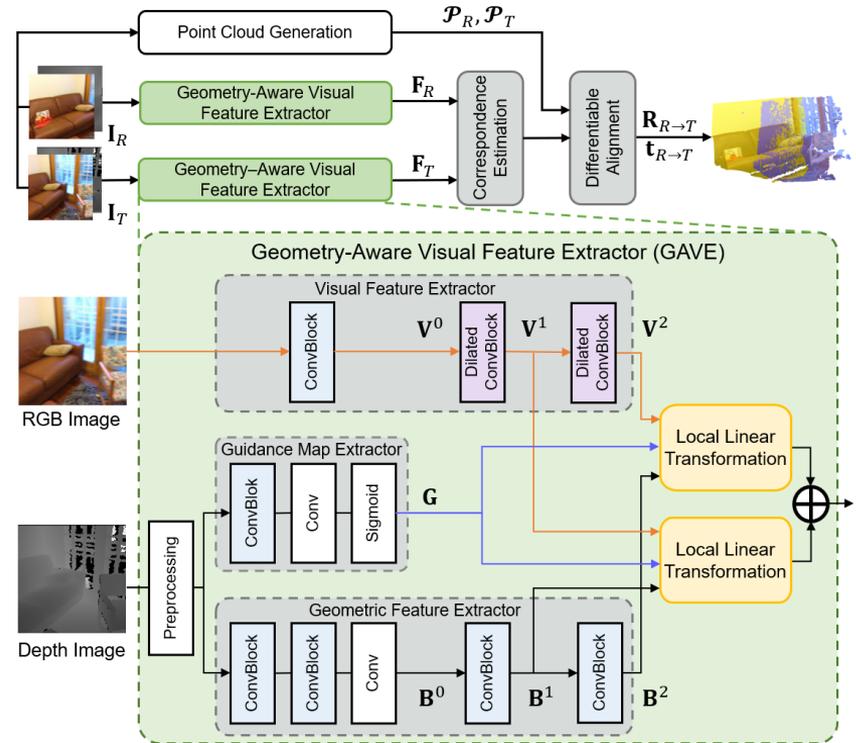


# Fazit und Zukunft



# Zukunft

- Weiteres LoRA/SAM-Training, um Score zu erhöhen
- evtl. Anpassen der Bounding Boxes für Punktwolken Generierung
- Mapping des segmentierten Roboters mit RGB Daten auf das Modell Objekt (ICP)
- Response mit relativer Position an die Hololens (SAM on Hololens bzw. externer Server)



# Lernerfahrungen

Jan Philipp	Benedikt
Tools: tensorflow 2 / Keras / Pytorch / Python	Tools: tensorflow 2 / Keras / Pytorch / Python / Blender
Paper: Segment Anything, LoRA	Paper: Segment Anything, LoRA
Einrichtung, Training und Abspielen von Neuronalen Netzen, Änderung von Layer und Parametern, Erstellung von Train/Test Datensätzen (für fine-tuning)	Erstellung von Train/Test Datensätzen (für fine-tuning) und Evaluierung, Training und Abspielen von Neuronalen Netzen

# Lernerfahrungen

<b>Christoph</b>	<b>Elmar</b>
Tensorflow und Pytorch, HL2SS (Research Mode), Jupyter, Open3D, CloudCompare	Tensorflow, Pytorch, Cloudcompare, HL2SS, Hololens 2 Datenextraktion
Paper: Swin3D [e], HL2SS [b], Joint 2D-3D-Semantic Data for Indoor Scene Understanding [d]	Paper: Swin3D [e], HL2SS [b], Uni3DScenes Structured3D [d]
Manuelle und automatische Punktwolkenverarbeitung, -registrierung und Labelling, Finetuning/Evaluation Segmentierung, Punktwolkenregistrierung (ICP)	Datensegmentierung (manuelles Labeling) Pre Trained Modell finetuning Transformer Anwendung

# Literaturverzeichnis

- [a] <https://github.com/facebookresearch/segment-anything>
- [g] <https://github.com/facebookresearch/segment-anything/issues/4>
- [h] <https://keras.io/examples/vision/sam/>
- [i] [https://keras.io/guides/keras\\_cv/segment\\_anything\\_in\\_keras\\_cv/](https://keras.io/guides/keras_cv/segment_anything_in_keras_cv/)
- [j] <https://github.com/JamesQFreeman/LoRA-ViT>
- [k] <https://arxiv.org/pdf/2311.08236.pdf>
- [l] <https://www.youtube.com/watch?v=DhRoTONcyZE>
- [b] Dibene, Juan Carlos und Dunn, Enrique: HoloLens 2 Sensor Streaming (2022) <https://arxiv.org/pdf/2211.02648.pdf>
- [c] Joint 2D-3D-Semantic Data for Indoor Scene Understanding  
[http://buildingparser.stanford.edu/images/2D-3D-S\\_2017.pdf](http://buildingparser.stanford.edu/images/2D-3D-S_2017.pdf)
- [d] Structured3D: A Large Photo-realistic Dataset for Structured 3D Modeling  
[https://www.ecva.net/papers/eccv\\_2020/papers\\_ECCV/papers/123540494.pdf](https://www.ecva.net/papers/eccv_2020/papers_ECCV/papers/123540494.pdf)
- [e] Swin3D: A Pretrained Transformer Backbone for 3D Indoor Scene Understanding  
<https://arxiv.org/abs/2304.06906>
- [f] Armeni, Iro et al.: 3D Semantic Parsing of Large-Scale Indoor Spaces – Supplementary Material (2016)  
[https://svl.stanford.edu/assets/publications/pdfs/3dsemanticparsing\\_supp\\_mat.pdf](https://svl.stanford.edu/assets/publications/pdfs/3dsemanticparsing_supp_mat.pdf)

## Tipps und Anmerkungen von Varanasi

- **Besser Zitieren (mit Autoren etc)**
- **Ist auch ok wenn wir mehr als 8 seiten Projektbericht haben**
- **das projekt soll nicht in der schublade verschwinden → neue Studenten akquirieren**
- **auch auf das automatische Labelling eingehen**
- **wir sollen auch auf die technischen details eingehen, zb auch bei der Erstellung des Datensatzes**
- **das sind quasi 3 Projekte **
- **wir sollen die Herausforderungen erwähnen, sodass die nach uns, da weiterarbeiten können**