**University of Brasília - UnB**

Institute of Exact Sciences
Department of Computer Science

# Towards Complete 3D Indoor Scene Understanding from a Single Point-of-View

Qualifying examination of the Ph.D. Program in Computer Science

**Aloisio Dourado Neto**

Supervisor
Prof. Dr. Teófilo Emidio de Campos

Comitee

Prof. Dr. Anderson de Rezende Rocha
Unicamp

Prof. Dr. Ricardo de Queiroz
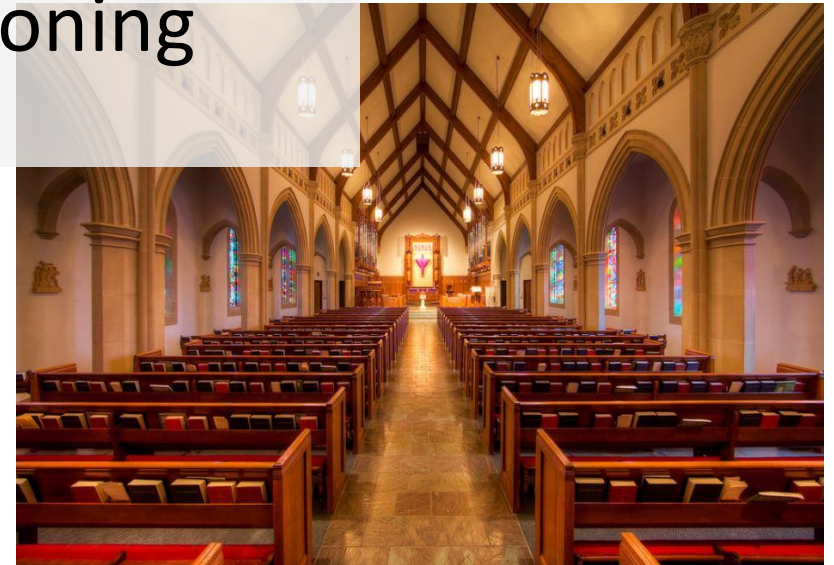CIC/UnB

# Presentation Outline

- Introduction
  - Motivation
  - Problem statement
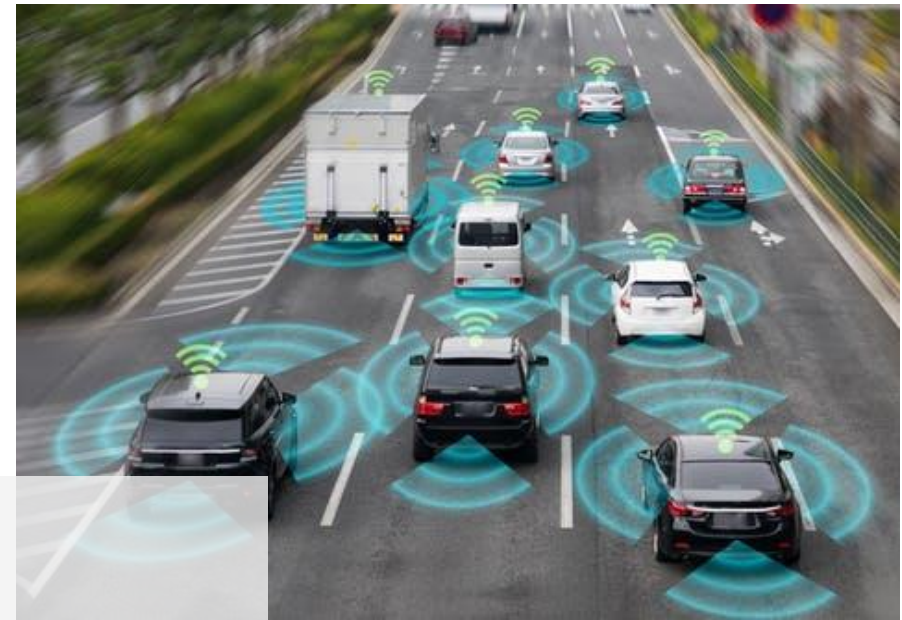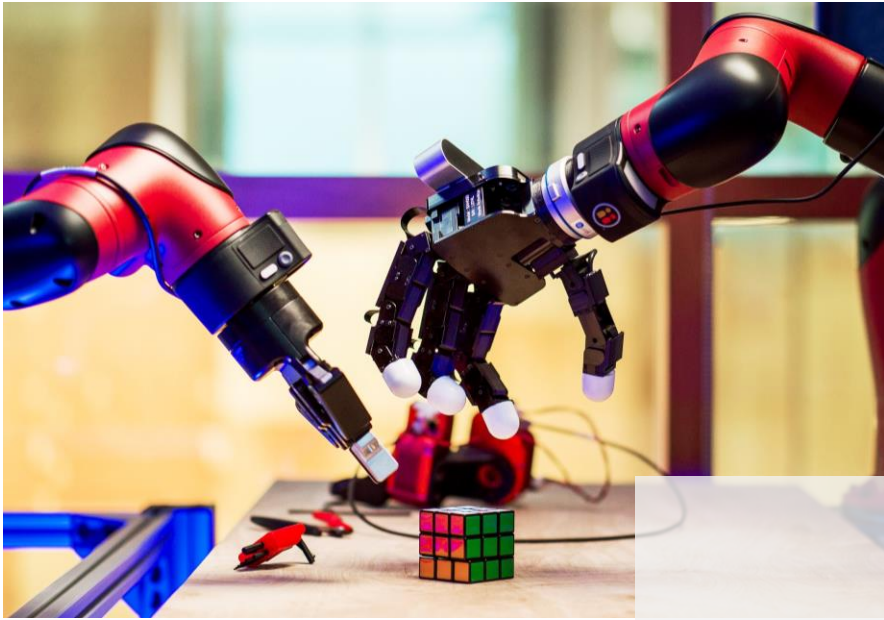  - Objectives
- Research steps
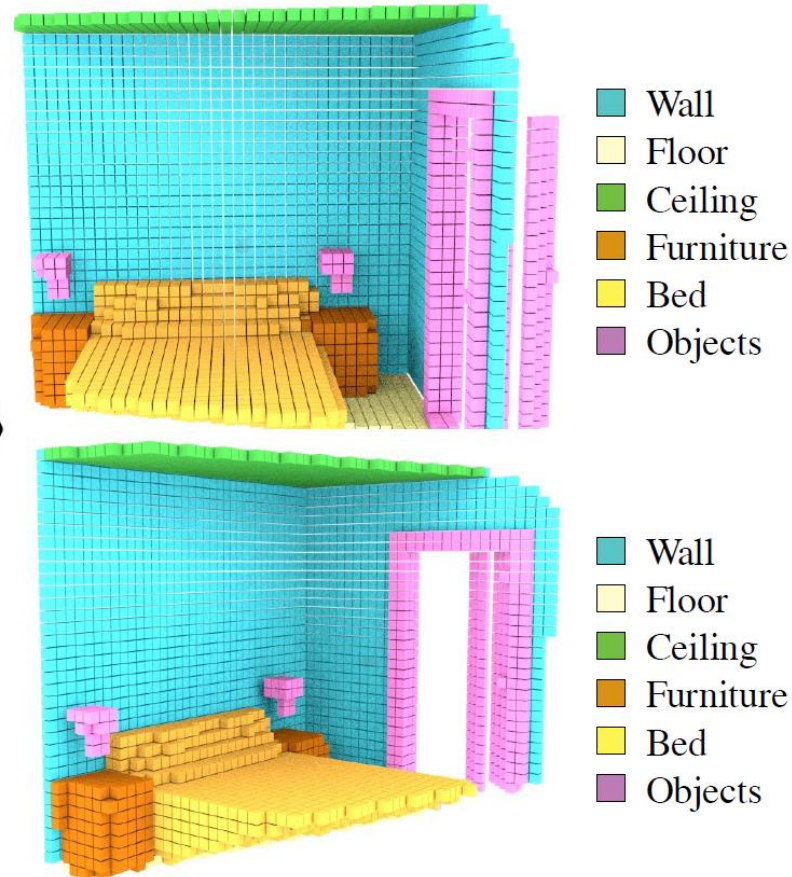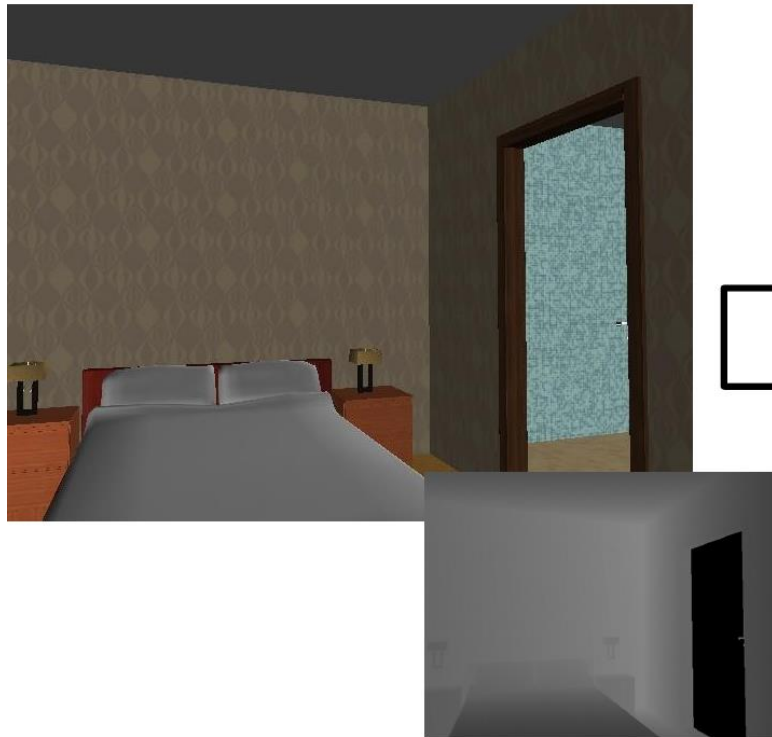- Work plan

Introduction



3D Scene Reasoning

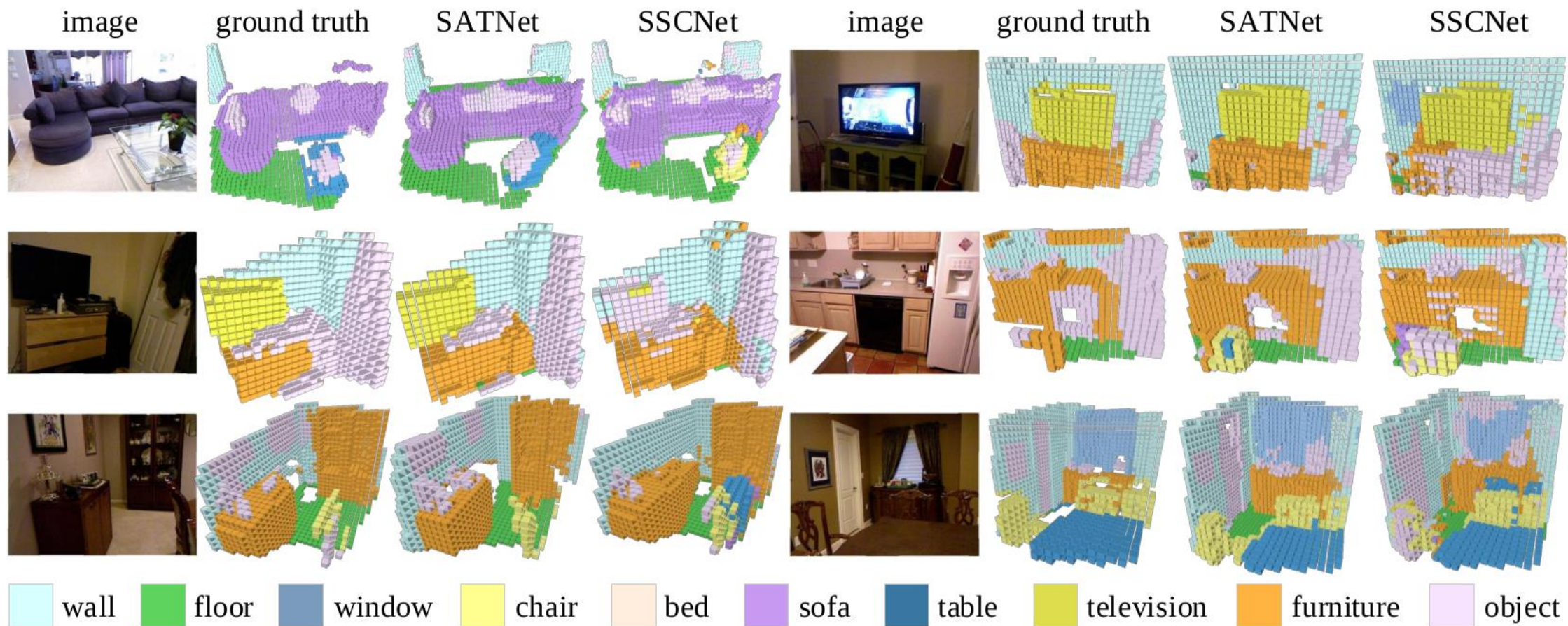# Applications

# Semantic Scene Completion



Introduced by Song *et al.*[107] in 2017

Trained a 3D CNN that jointly deals with both completion and semantic segmentation

[107] Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., and Funkhouser, T.: Semantic Scene Completion from a Single Depth Image. In Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, Hawaii, July 21-26, pp. 190–198, Piscataway, NJ, July 2017. IEEE. 2, 3, 4, 18, 45, 46, 47, 51, 52, 53, 64, 68, 70
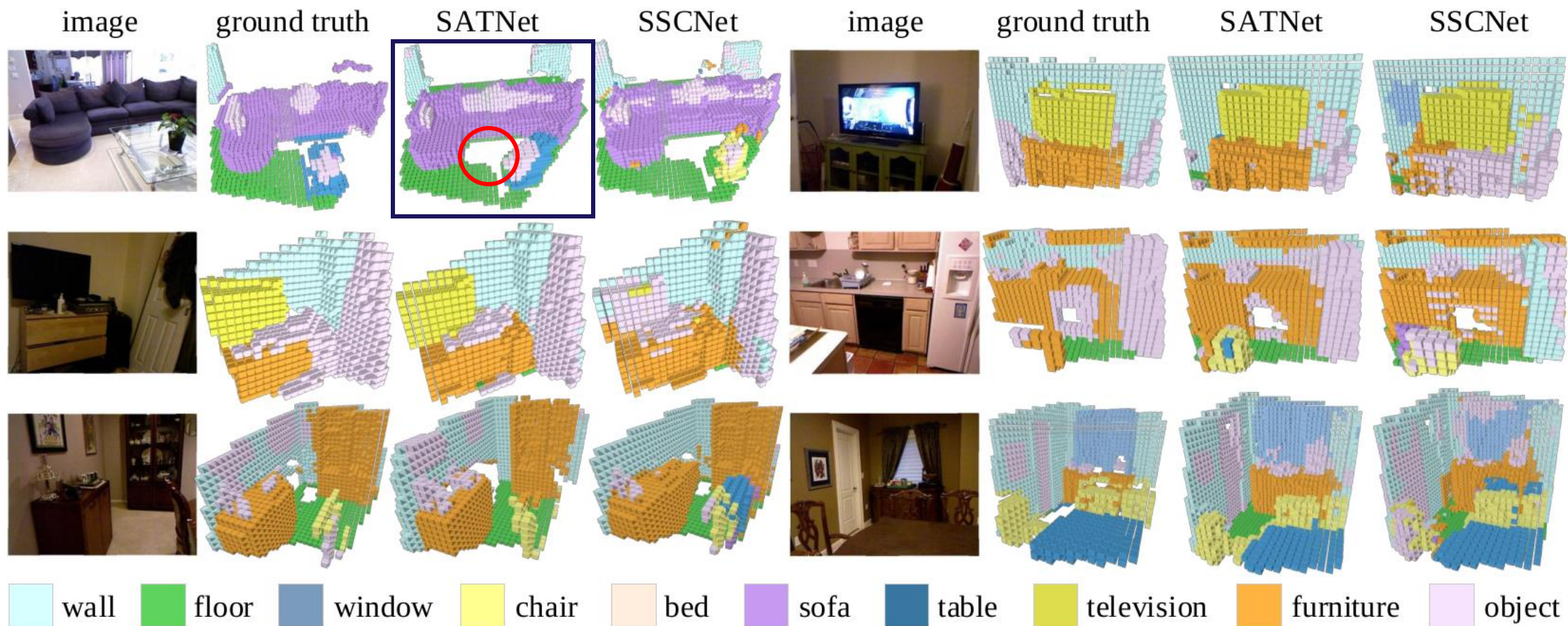
# Problem Statement



Qualitative results on NYUv2 dataset from Liu *et al.* [70]

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59
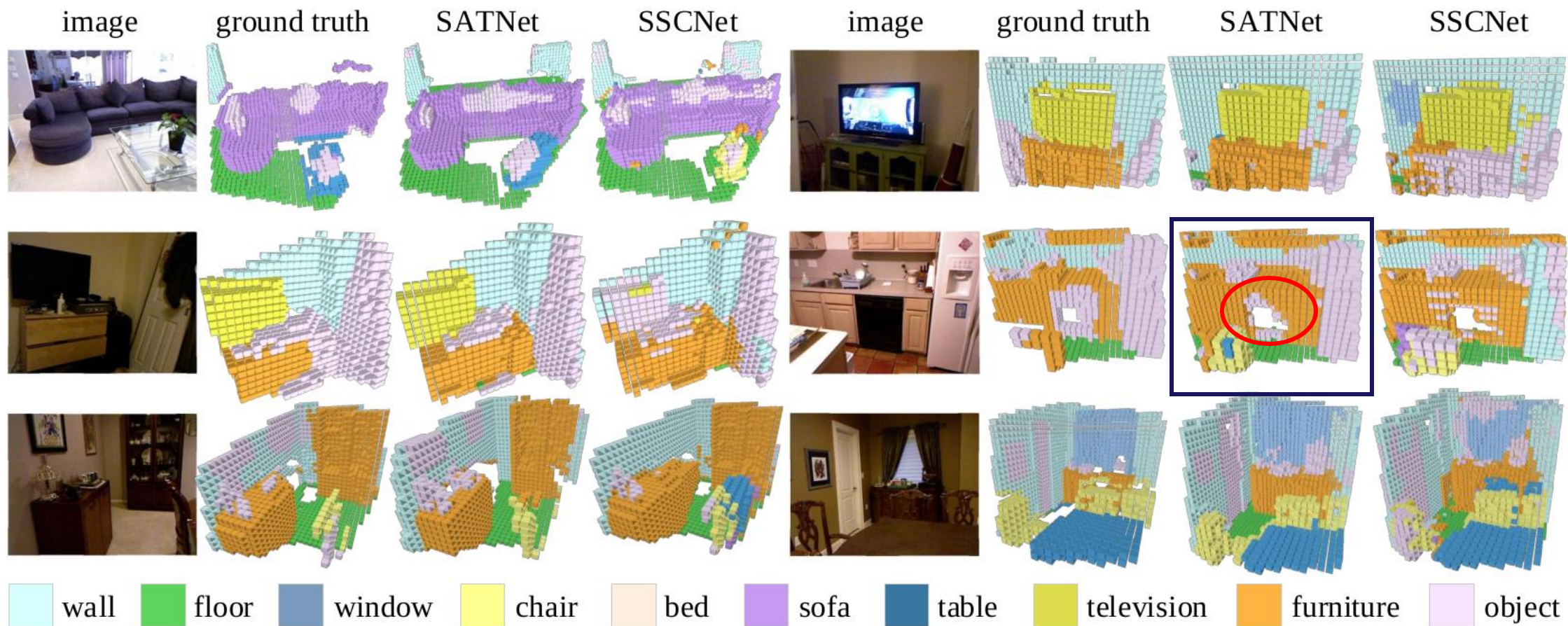
# Problem Statement



Qualitative results on NYUv2 dataset from Liu *et al.* [70]

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59
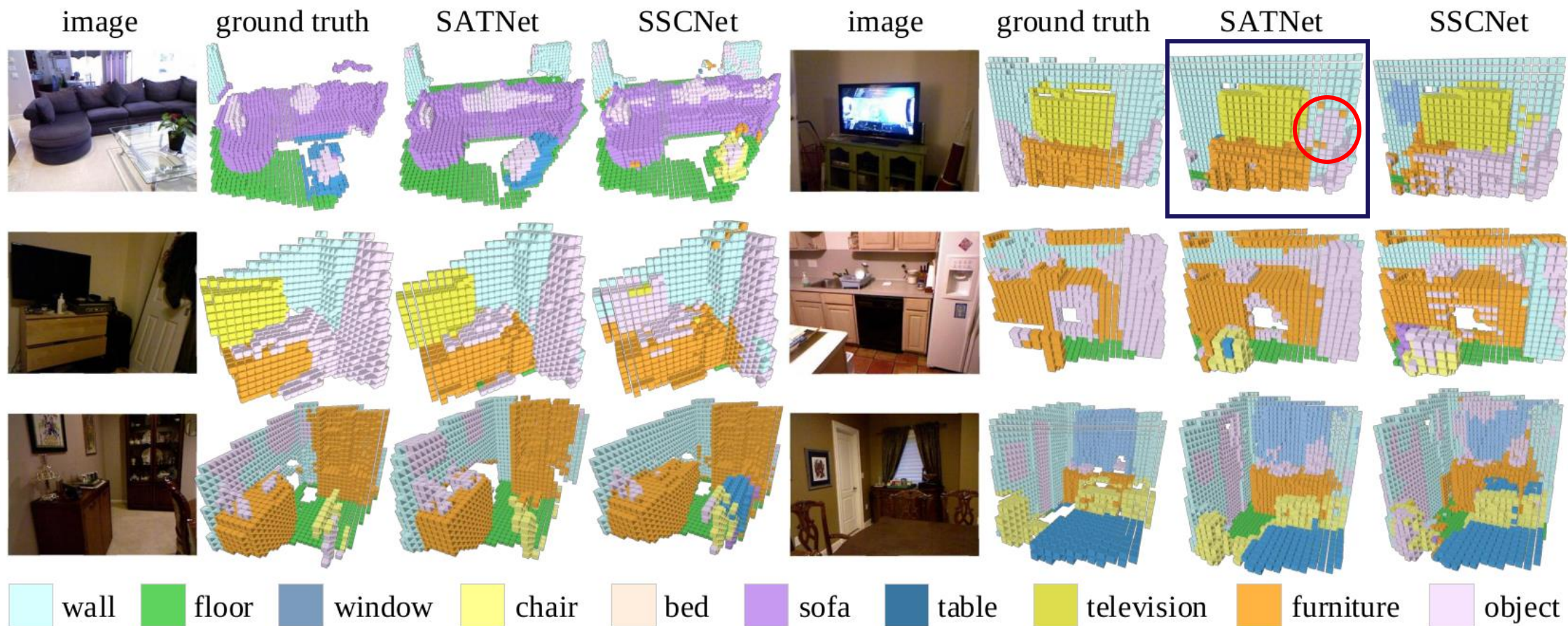
# Problem Statement



Qualitative results on NYUv2 dataset from Liu *et al.* [70]

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59
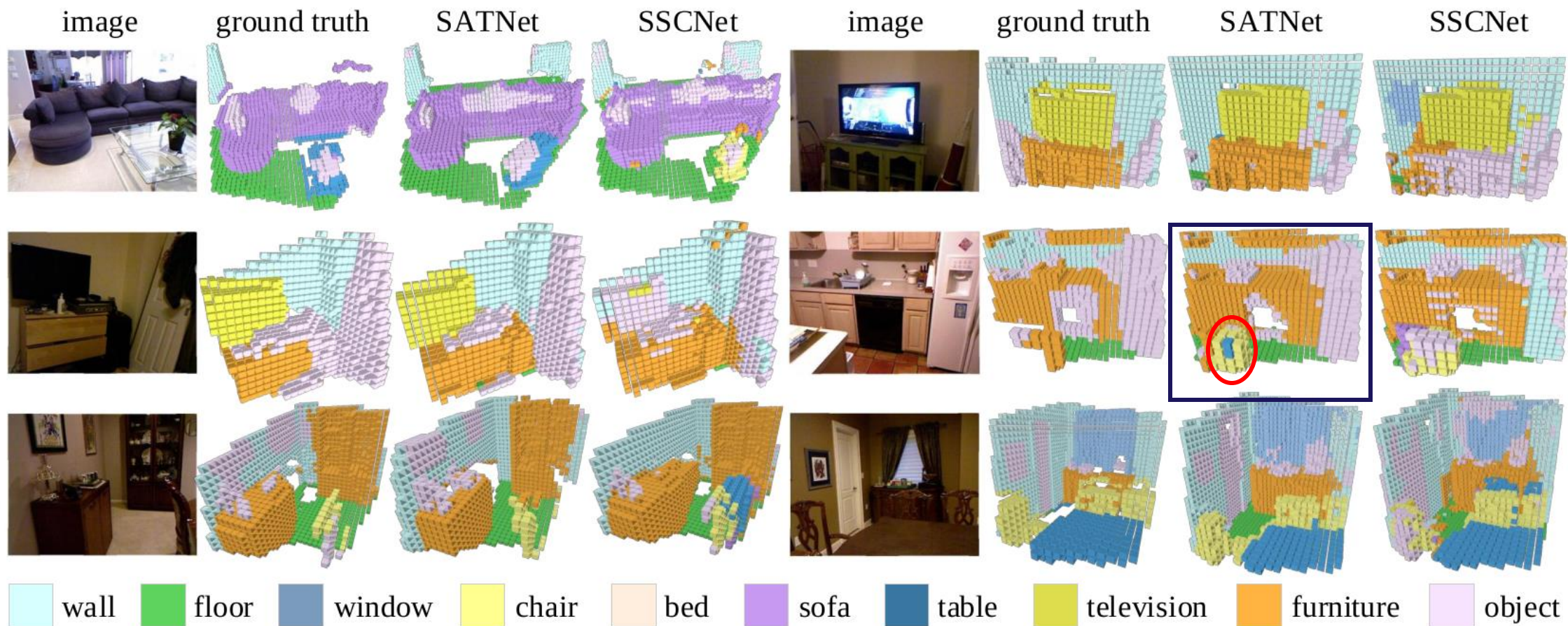
# Problem Statement



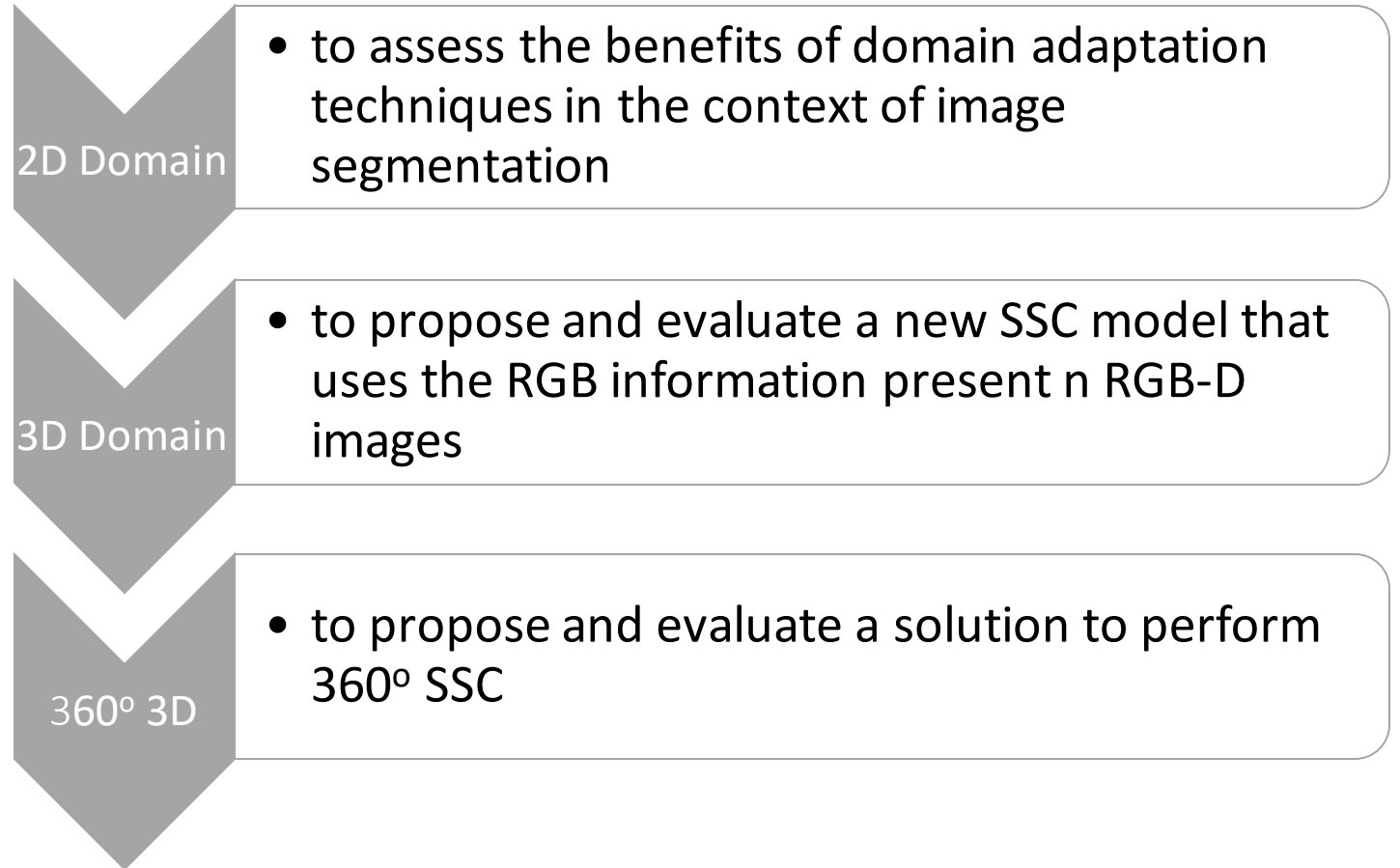Qualitative results on NYUv2 dataset from Liu *et al.* [70]

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59

# Problem Statement



Qualitative results on NYUv2 dataset from Liu *et al.* [70]

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59

# Problem Statement

- Two main deficiencies of current approaches:
  - the RGB part of the RGB-D image is not completely explored;
  - they are limited to the restricted FOV of depth sensors like Kinect

# Objectives

New tools and models that could push SSC solutions towards a complete understating of the whole indoor scene

**2D Domain**
- to assess the benefits of domain adaptation techniques in the context of image segmentation

**3D Domain**
- to propose and evaluate a new SSC model that uses the RGB information present n RGB-D images

**360º 3D**
- to propose and evaluate a solution to perform 360º SSC

# Chapter 3

## Fully Convolutional Networks, Domain Adaptation and Semantic Segmentation

# Fully Convolutional Networks, Domain Adaptation and Semantic Segmentation

## Why work on 2D?

- Work on 3D is hard
- Less previous works to compare!
- Start to explore domain adaptation and segmentation in an easier domain

[53] Kakumanu, P., Makrogiannis, S., and Bourbakis, N.: A survey of skin-color modeling and detection methods. Pattern Recognition, 40(3):1106 − 1122, 2007, ISSN 0031-3203. 27

[12] Brancati, N., Pietro, G.D., Frucci, M., and Gallo, L.: Human skin detection through correlation rules between the YCb and YCr subspaces based on dynamic color clustering. Computer Vision and Image Understanding, 155:33 − 42, 2017, ISSN 1077-3142. 27, 28, 35, 36, 39, 42

# Fully Convolutional Networks, Domain Adaptation and Semantic Segmentation

- Why the skin segmentation application?
  - Research field where some criticisms regarding the use of CNNs/FCNs are made:
    - the need for large training datasets [53]
    - the specificity or lack of generalization of neural nets
    - long prediction time [12]
  - We wanted to try to refute those criticisms

[53] Kakumanu, P., Makrogiannis, S., and Bourbakis, N.: A survey of skin-color modeling and detection methods. Pattern Recognition, 40(3):1106 − 1122, 2007, ISSN 0031-3203. 27

[12] Brancati, N., Pietro, G.D., Frucci, M., and Gallo, L.: Human skin detection through correlation rules between the YCb and YCr subspaces based on dynamic color clustering. Computer Vision and Image Understanding, 155:33 − 42, 2017, ISSN 1077-3142. 27, 28, 35, 36, 39, 42

# Previous Works

Historically, color-based or texture methods were preferred [49, 100]

Current state-of the-art works still rely on local approaches:

- Skin-color separation [12, 33]
- Patch-based CNN [74]

The use of domain adaptation methods for this problem is not common

[12] Brancati, N., Pietro, G.D., Frucci, M., and Gallo, L.: Human skin detection through correlation rules between the YCb and YCr subspaces based on dynamic color clustering. Computer Vision and Image Understanding, 155:33 − 42, 2017, ISSN 1077-3142.
27, 28, 35, 36, 39, 42
[33] Faria, R.A.D. and Hirata Jr., R.: Combined correlation rules to detect skin based on dynamic color clustering. In Proceedings of the 13th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISAPP), vol. 5, pp. 309–316. INSTICC, SciTePress, 2018, ISBN 978-989-758-290-5. 28, 35, 36
[49] Huynh-Thu, Q., Meguro, M., and Kaneko, M.: Skin-Color-Based Image Segmentation and Its Application in Face Detection. In MVA, pp. 48–51, 2002. 27, 39
[74] Lumini, A. and Nanni, L.: Fair comparison of skin detection approaches on publicly available datasets. Techn. rep., Cornell University Library, CoRR/cs.CV, August 2019. arXiv:1802.02531 (v3). 28, 43
[100] Shrivastava, V.K., Londhe, N.D., Sonawane, R.S., and Suri, J.S.: Computer-aided diagnosis of psoriasis skin images with HOS, texture and color features. Comput. Methods Prog. Biomed., 126(C):98–109, Apr. 2016, ISSN 0169-2607. 27

# Experiments

## In-domain:

- Local CNN vs Holistic FCN
- Comparison to current color-based state-of-the-art

## Cross-domain:

- Assessment of 3 simple methods



Fine-tuning

Pseudo-label

Combined approach

# Models

## Local, Patch-based CNN

35 x 35 x 3

35 x 35 x 16

17 x 17 x 32

8 x 8 x 32

4 x 4 x 16

32

1

- Input
- 2D Convolution(3x3) + ReLu + MaxPooling
- Flatten
- Dense
- Dense + Sigmoid

## Holistic, u-shaped FCN

768 x 768 x 16

MaxPooling(2x2) Strides(2x2)

UpSampling2D

384 x 384 x 32

MaxPooling(2x2) Strides(2x2)

UpSampling2D

192 x 192 x 64

MaxPooling(2x2) Strides(2x2)

UpSampling2D

96 x 96 x 128

UpSampling2D

MaxPooling(2x2) Strides(2x2)

- Conv2D(3x3) + ReLU + BatchNorm
- Concatenate
- Conv2D(1x1) + Sigmoid

# Datasets

SFA[15]
(1,118 images)



Pratheepan[117]
(78 images)

Compaq[51]
(4,670 images)

VPU[93]
(290 images)

# Supervised Training vs Domain Adaptation



Comparison of source only vs. domain adaptation combined approach in the Compaq→Pratheepan scenario

# Conclusions

Refuted criticisms regarding the use of Deep Convolutional Networks for skin segmentation

- Color or texture separation may suffice:
  - Our two CNN approaches performed much better than the color-based state-of-the-art

- CNNs are slow:
  - Our U-Net inference time was enough for real-time applications

- CNNs need too much data to generalize:
  - With no labeled data -> 60% improvement

# Publication

*Domain Adaptation for Holistic Skin Detection*



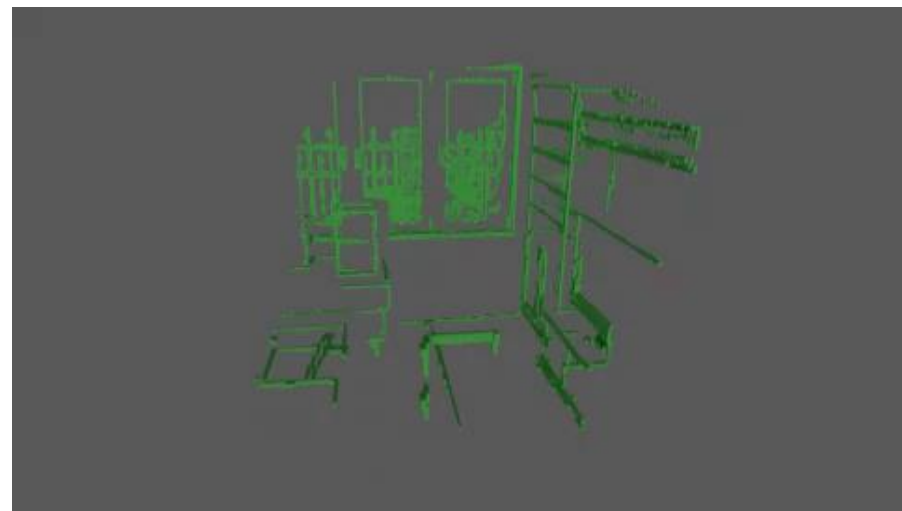*Submitted to International Journal of Pattern Recognition and Artificial Intelligence (Capes Qualis B1)

[30] Dourado, A., Guth, F., de Campos, T.E., and Weigang, L.: Domain adaptation for holistic skin detection. Tech. Rep. arXiv:1903.0969, Cornell University Library, 2019. http://arxiv.org/abs/1903.06969. 6, 26

# Using RGB Edges to improve Semantic Scene Completion from RGB-D Images

## Chapter 4

# Previous Works

## Depth maps only

- **SSCNET: Song et al. [107]**
  - Seminal paper
  - Proposed F-TSDF encoding
  - Introduced SUNCG Dataset

[107] Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., and Funkhouser, T.: Semantic Scene Completion from a Single Depth Image. In Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, Hawaii, July 21-26, pp. 190–198, Piscataway, NJ, July 2017. IEEE. 2, 3, 4, 18, 45, 46, 47, 51, 52, 53, 64, 68, 70

# Previous Works

## Depth maps only

- Guo and Tong [40]:
  - 2D features projected to 3D



(a) SSCNet Structure
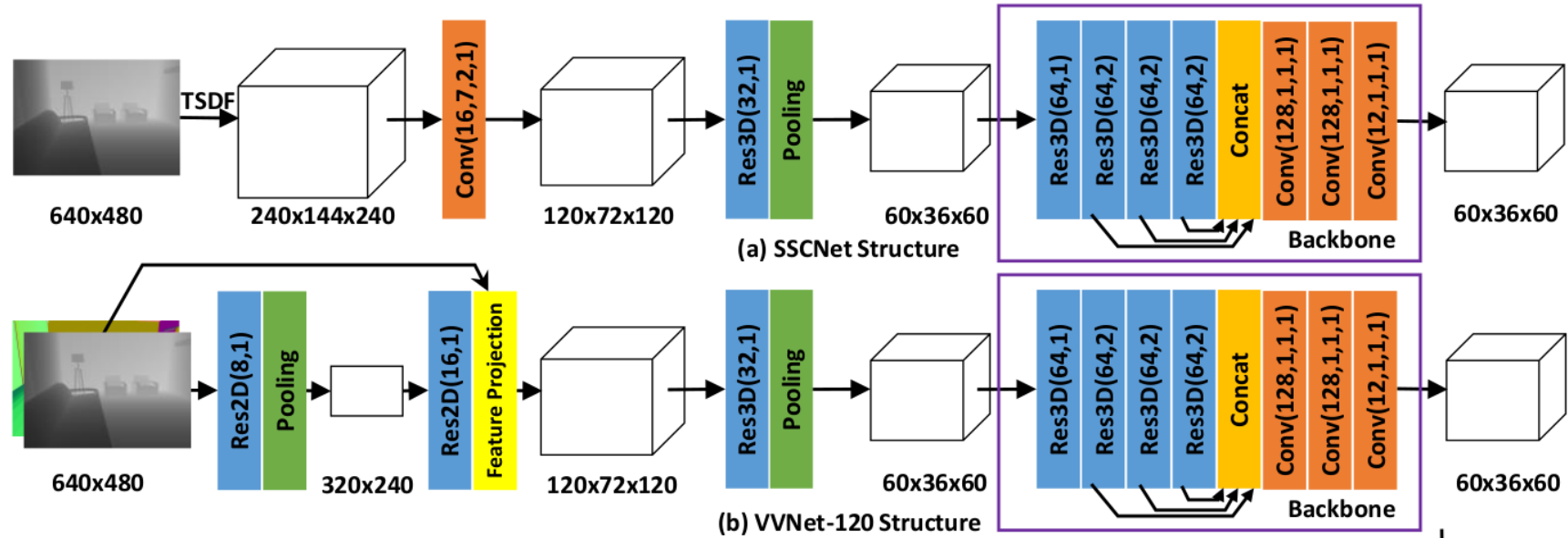
(b) VVNet-120 Structure

[40] Guo, Y. and Tong, X.: View-Volume Network for Semantic Scene Completion from a Single Depth Image. In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, pp. 726–732, Stockholm, Sweden, July 2018. International Joint Conferences on Artificial Intelligence Organization, ISBN 978-0-9992411-2-7. https://doi.org/10.24963/ijcai.2018/101. 2, 4, 18, 46, 52, 53

# Previous Works

## Depth maps only

- Guo and Tong [40]:
  - 2D features projected to 3D



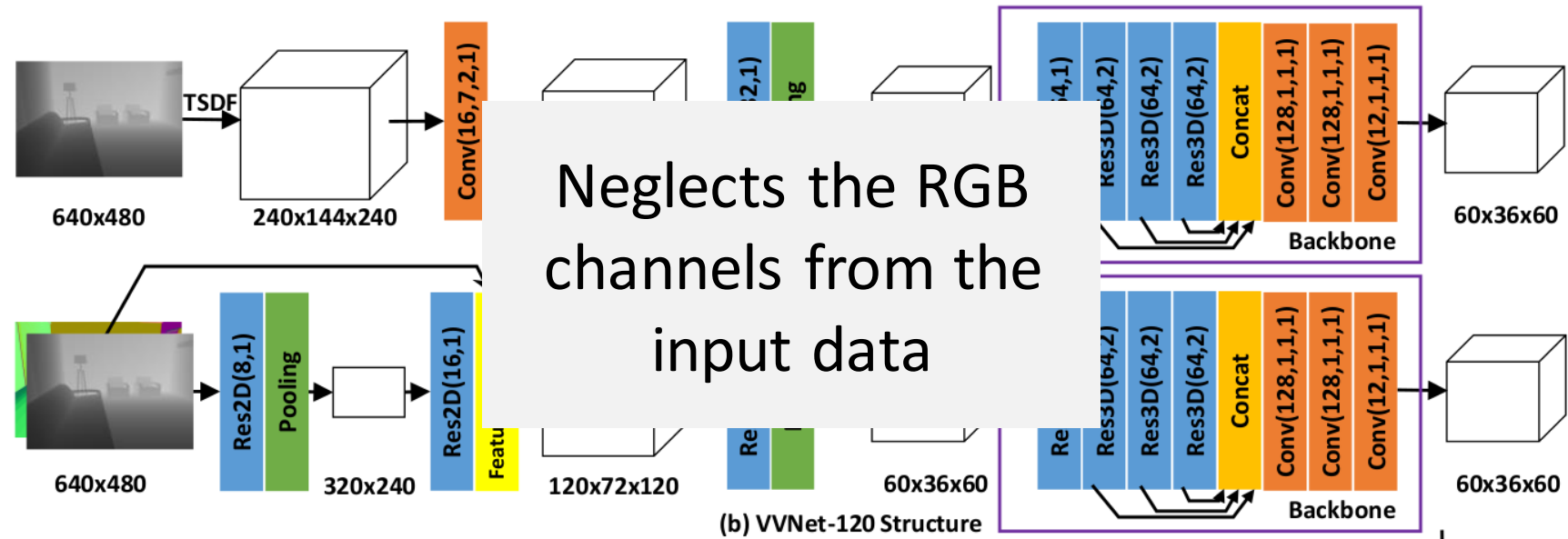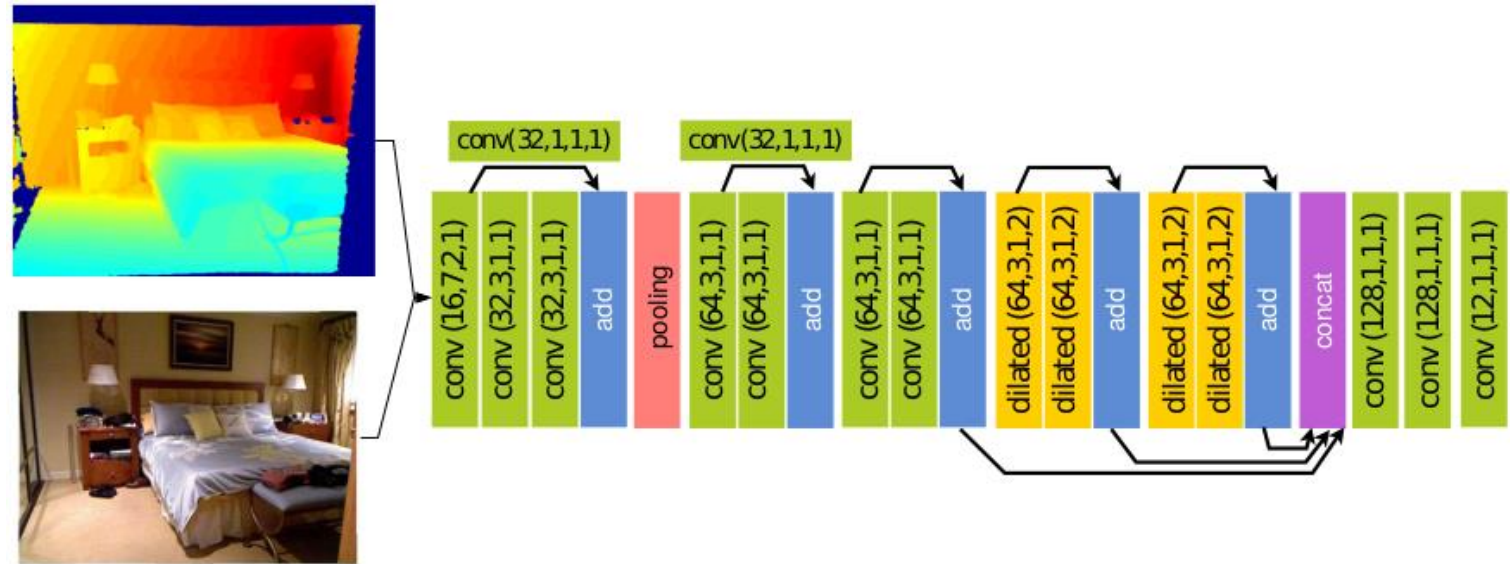Neglects the RGB channels from the input data

[40] Guo, Y. and Tong, X.: View-Volume Network for Semantic Scene Completion from a Single Depth Image. In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, pp. 726–732, Stockholm, Sweden, July 2018. International Joint Conferences on Artificial Intelligence Organization, ISBN 978-0-9992411-2-7. https://doi.org/10.24963/ijcai.2018/101. 2, 4, 18, 46, 52, 53

# Previous Works

## Depth maps plus RGB

- Guedes *et al.*[38]

[38] Guedes, A.B.S., de Campos, T.E., and Hilton, A.: Semantic scene completion combining colour and depth: preliminary experiments. In ICCV workshop on 3D Reconstruction Meets Semantics (3DRMS), Venice, Italy, October 2017.
Event webpage: http://trimbot2020.webhosting.rug.nl/events/events-2017/3drms/. Also published at arXiv:1802.04735. 4, 45, 46, 47, 52, 53

# Previous Works

## Depth maps plus RGB

- Guedes *et al.*[38]



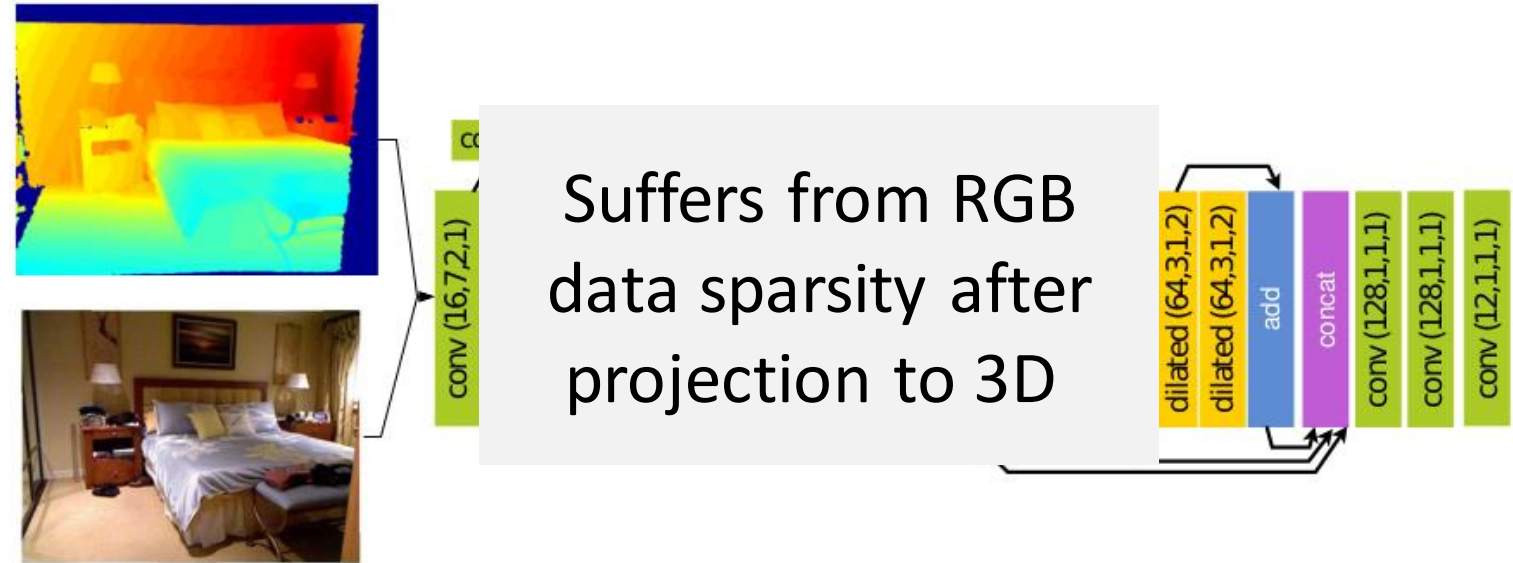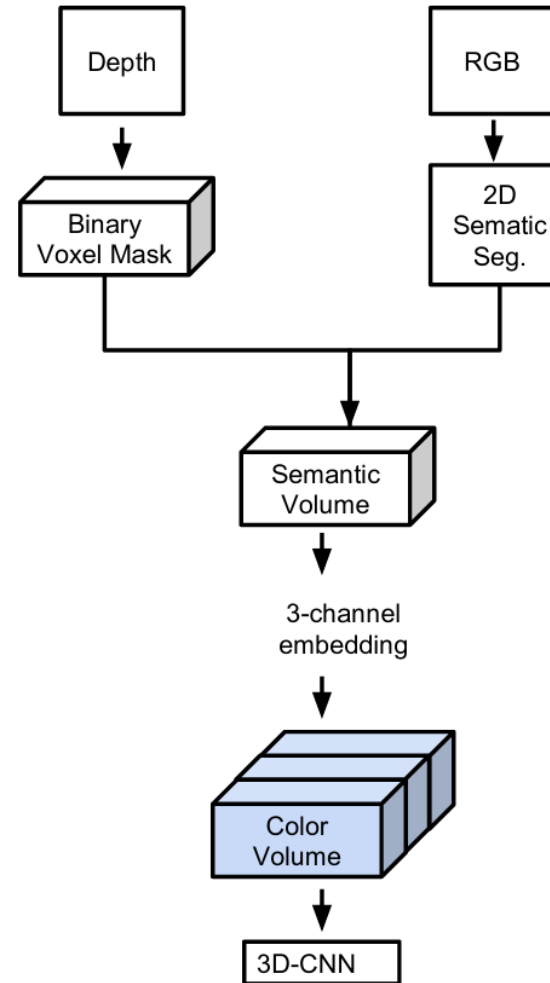Suffers from RGB data sparsity after projection to 3D

[38] Guedes, A.B.S., de Campos, T.E., and Hilton, A.: Semantic scene completion combining colour and depth: preliminary experiments. In ICCV workshop on 3D Reconstruction Meets Semantics (3DRMS), Venice, Italy, October 2017. Event webpage: http://trimbot2020.webhosting.rug.nl/events/events-2017/3drms/. Also published at arXiv:1802.04735. 4, 45, 46, 47, 52, 53

# Previous Works

## Depth map plus 2D segmentation

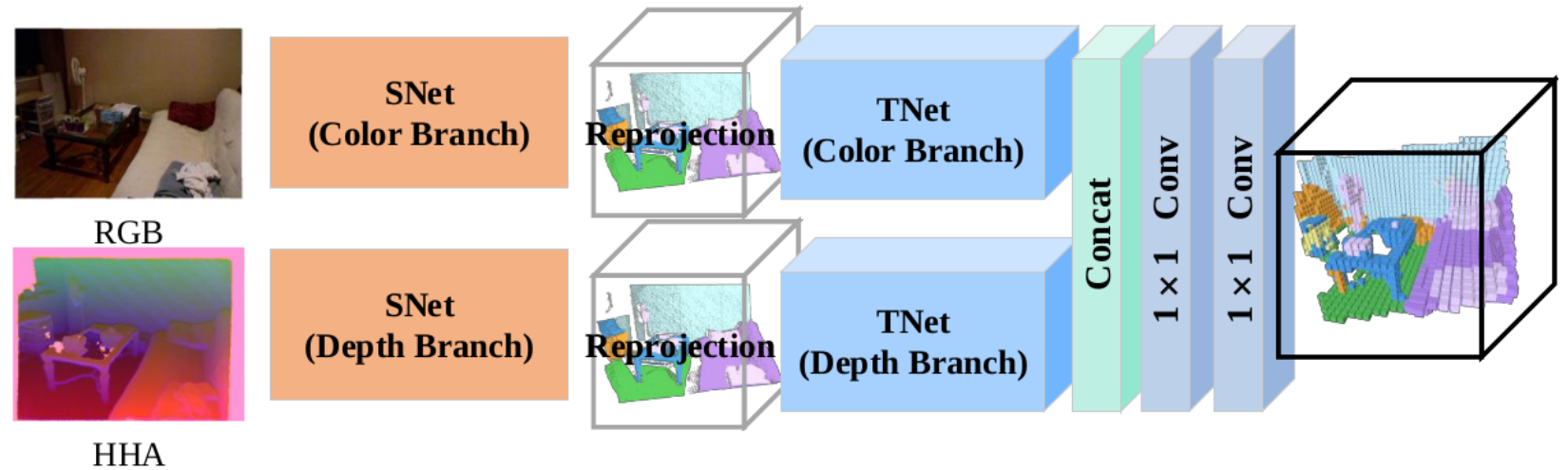- Two stream 3D semantic scene completion: Garbade *et al.*[36]



[36] Garbade, M., Sawatzky, J., Richard, A., and Gall, J.: Two stream 3D semantic scene completion. Tech. Rep. arXiv:1804.03550, Cornell University Library, 2018. http://arxiv.org/abs/1804.03550. 4, 45, 47, 52, 53

# Previous Works

## Depth map plus 2D segmentation

- TNetFusion: Liu *et al.* [70]



RGB

HHA

SNet (Color Branch) — Reprojection — TNet (Color Branch)

SNet (Depth Branch) — Reprojection — TNet (Depth Branch)
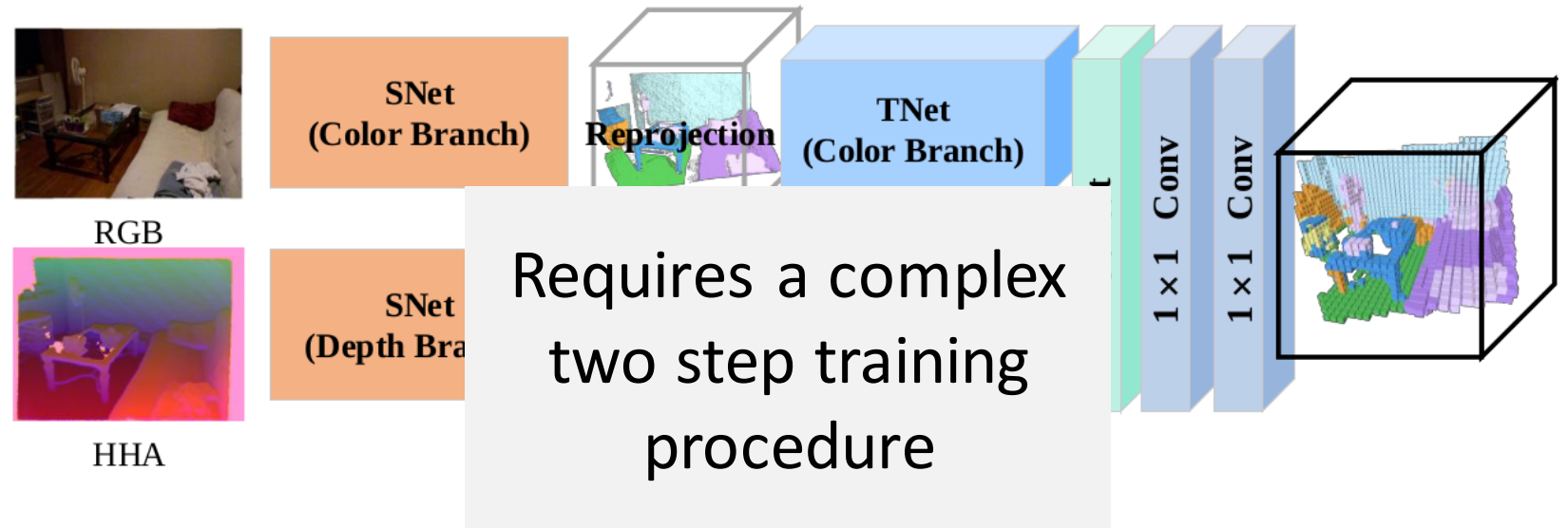
Concat — 1 × 1 Conv — 1 × 1 Conv

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59
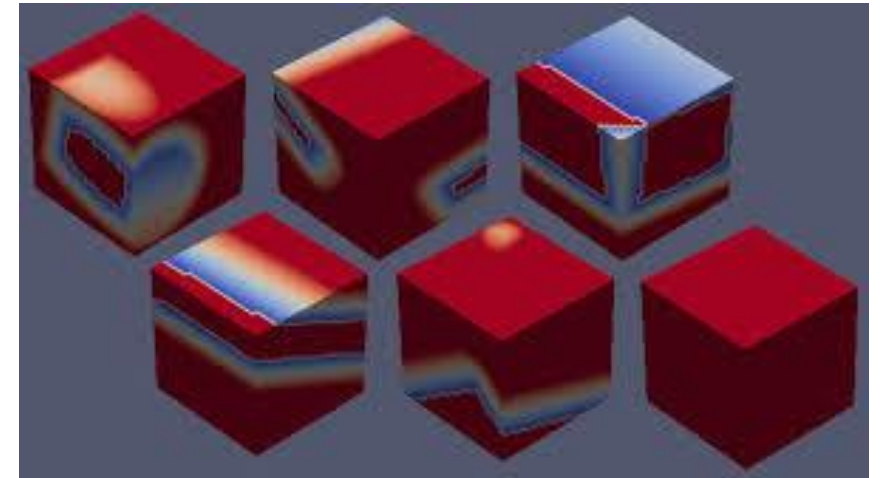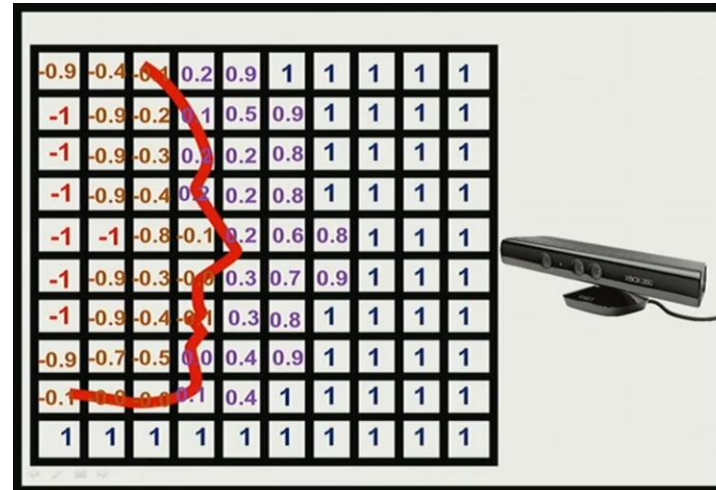
# Previous Works

## Depth map plus 2D segmentation

- TNetFusion: Liu *et al.*[70]



Requires a complex two step training procedure

[70] Liu, S., HU, Y., Zeng, Y., Tang, Q., Jin, B., Han, Y., and Li, X.: See and think: Disentangling semantic scene completion. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.): Procedings of Conference on Neural Information Processing Systems 31 (NIPS), pp. 263–274, Reed Hook, NY, 2018. Curran Associates, Inc. http://papers.nips.cc/paper/7310-see-and-think-disentangling-semantic-scene-completion. 2, 4, 45, 47, 52, 53, 58, 59

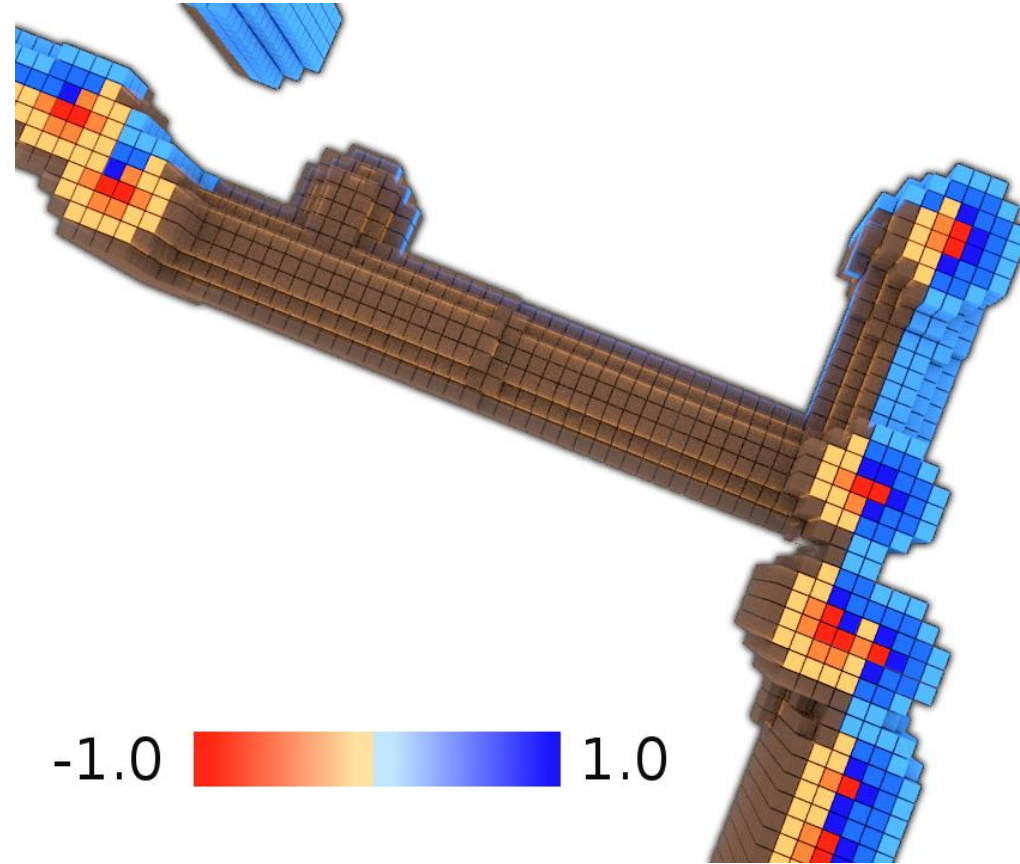# TSDF vs F-TSDF

- TSDF: Truncated Signed Distance Function



TSDF

## TSDF vs F-TSDF

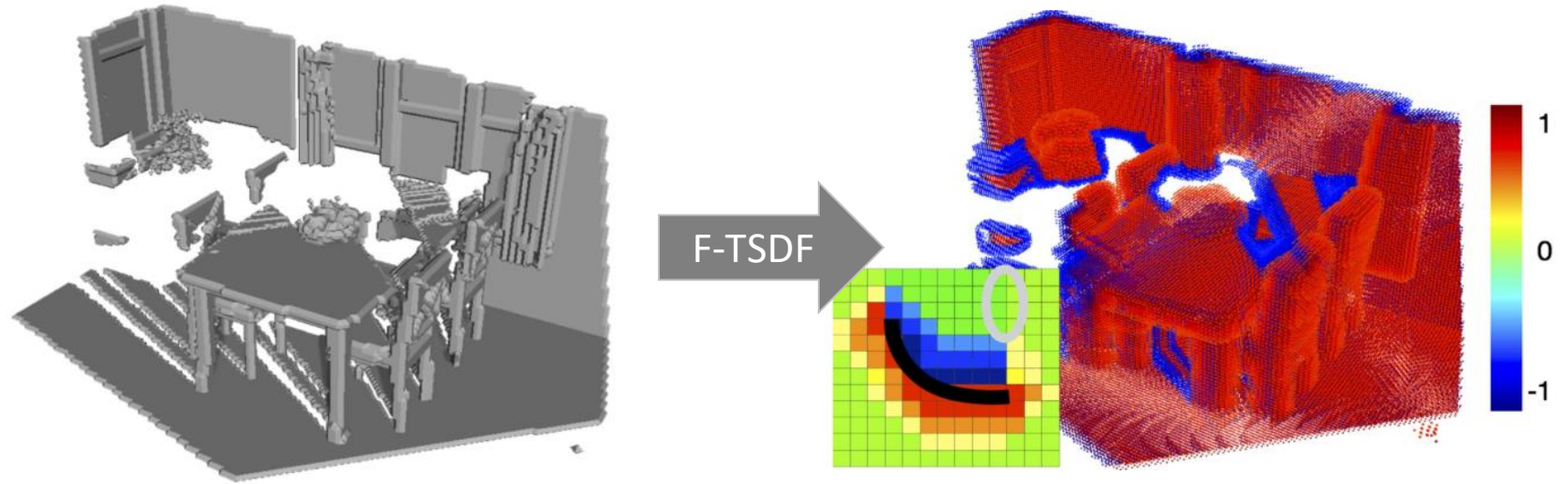- F-TSDF: Flipped Truncated Signed Distance Function



-1.0     1.0

F-TSDF

$$\text{F-TSDF} = \text{sign}(\text{TSDF}) \cdot (1-|\text{TSDF}|)$$
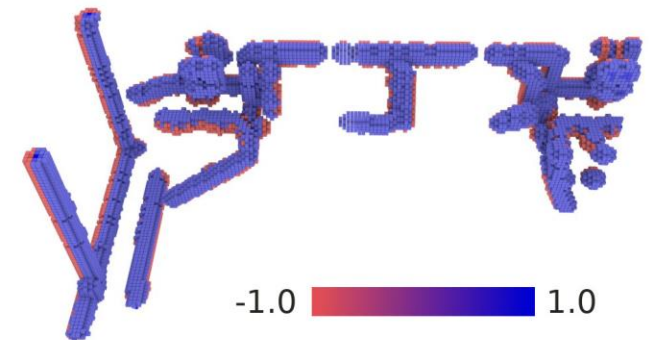
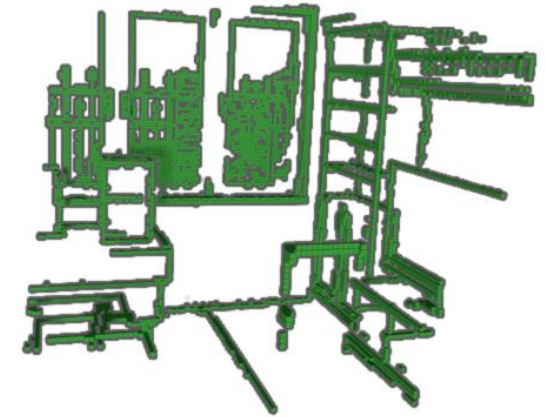## F-TSDF and the RGB Volume

- It is possible to apply F-TSDF to the occupancy volume



- However, RGB data is not binary!
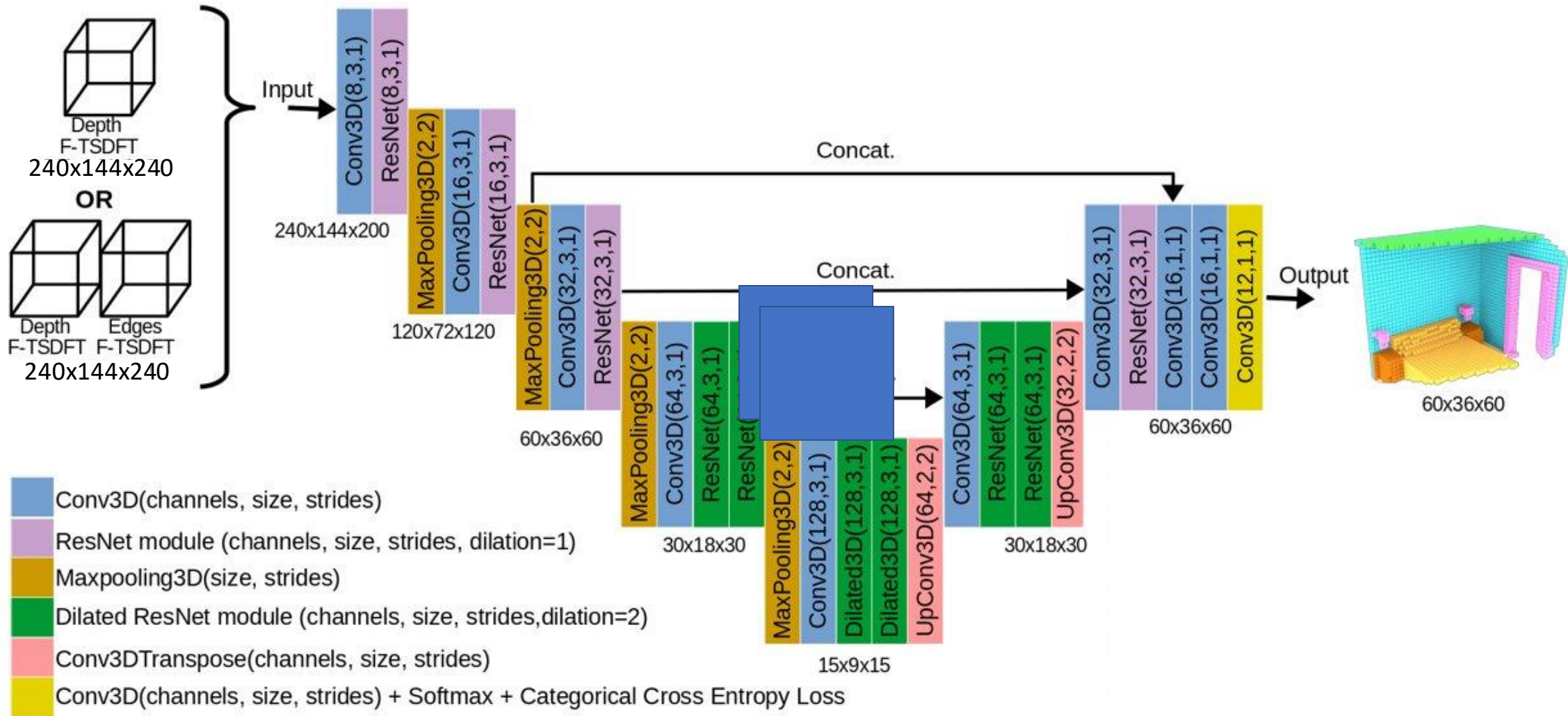
# Our Approach: EdgeNet

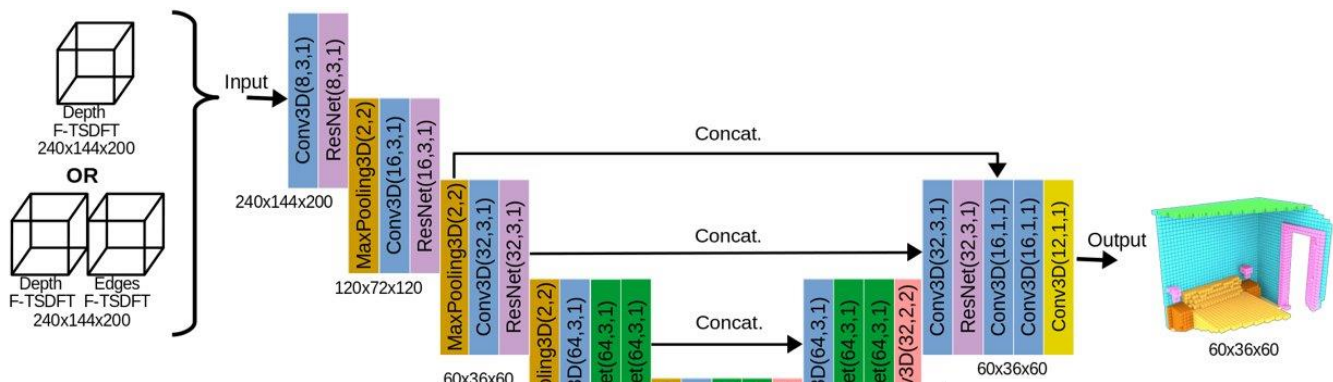- We extract information from RGB data using Canny Edge detector before F-TSDF

## Our implementation

- Offline F-TSDF calculation using portable C++ CUDA code

- We provide a software interface between CUDA and Python

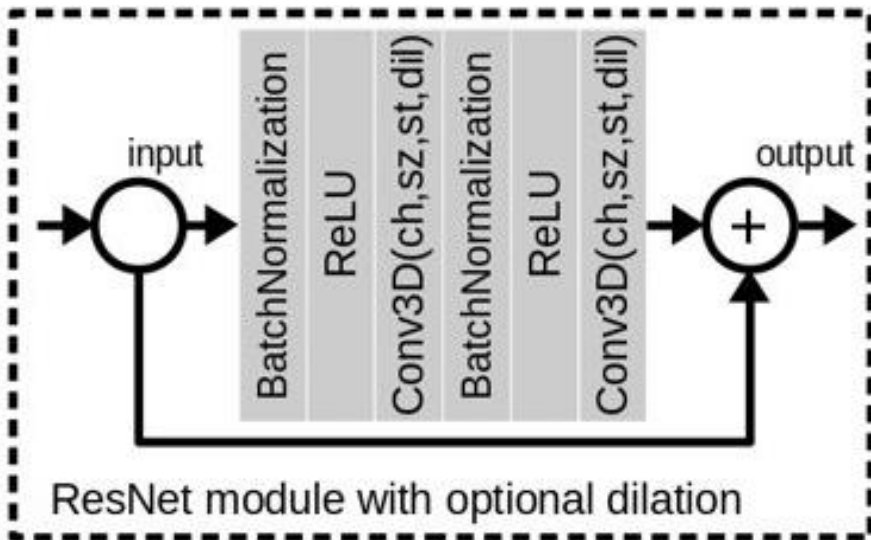- Preprocessing code is independent from the deep learning framework
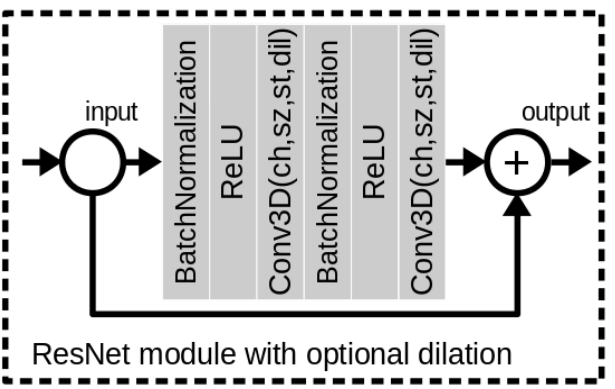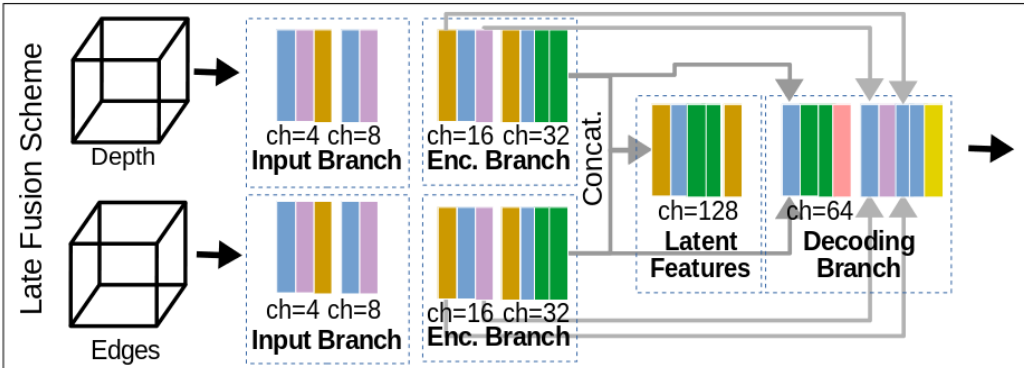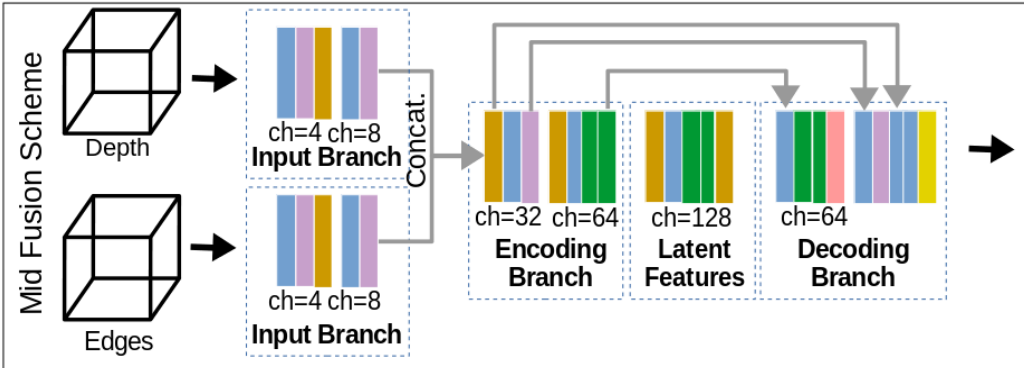
# Network Architecture

# Network Architecture

# Network Architecture - Fusion Schemes

# Network Architecture - Fusion Schemes

# Network Architecture - Fusion Schemes

# Network Architecture - Fusion Schemes

# Datasets

- ## SUNCG*



(a) SUNCG dataset     (b) 3D Scene     (c) Synthetic depth and volumetric ground truth

- ## NYUDv2**

*Song *et al*.[107]

**Silberman *et al*.[102]

# Training Time

- Ours
  - SUNCG: 4 days
  - NYU: 6 hours
- SSCNET
  - SUNCG: 7 days
  - NYU: 30 hours

# Quantitative Results

- New state-of-the-art result on SUNCG
- All new aspects of our solution contributed to the improvement
- Middle Fusion and Late Fusion schemes presented similar results on SUNCG
- Middle Fusion presented better results on NYUDV2

# Qualitative Results



Ground Truth            SSCNet            EdgeNet-MF

# Qualitative Results



Ground Truth                    SSCNet                    EdgeNet-MF

Higher overall accuracy

# Qualitative Results



Ground Truth                    SSCNet                    EdgeNet-MF

Hard-to-detect classes

# Qualitative Results



Ground Truth

SSCNet

EdgeNet-MF

NYU Ground Truth errors

# Conclusions

- A new end-to-end network architecture

- A new RGB enconding strategy

- Visually perceptible improvements

- Improvement over the state-of-the-art result on SUNCG

- We surpased  other end-to-end approaches on NYUv2

- An efficient and lightweight training pipeline for the task

# Publication

*EdgeNet: Sematic Scene Completion from a Single RGB-D Image*



*Accepted for publication in the proceedings of the 25th International Conference on Pattern Recognition (ICPR2020) (Capes Qualis A2)

[29] Dourado, A., de Campos, T.E., Kim, H., and Hilton, A.: EdgeNet: Semantic scene completion from RGB-D images. Tech. Rep. arXiv:1908.02893, Cornell University Library, 2019. http://arxiv.org/abs/1908.02893. 6, 44, 68

Extending
Semantic Scene
Completion for
$360^O$ Coverage

# Current Semantic Scene Completion Limitations



Regular RGB-D Sensor

Panoramic Image from
Matterport Camera

# Our approach



The 3DCNN is trained using SUNCG and fine-tuned in NYUDV2

This approach allows to use existing large and diverse RGB-D datasets for training.

# Results on Stanford 2D-3DS Dataset

| RGB Image | Input Volume | Predicted Volume | GT |
| --- | --- | --- | --- |



floor   wall   window   chair   table   sofa   furn.   objects

## Experiments on Spherical Stereo Images

- Stereo capture using commercial 360$^O$ cameras is one realistic approach to 360$^O$ SSC

- faster compared to Matterport scanning

- depth estimation is subject to errors due to occlusions between two camera views and correspondence matching errors

## Our approach

- vertical stereo setup
- Dense stereo matching with spherical stereo geometry [56]
- Depth map enhancement procedure:
  - Align the scene (Manhattan principle)
  - Apply Canny Edge Detector
  - RANSAC to fit a plane over coherent regions with similar colors



[56] Kim, H. and Hilton, A.: Block world reconstruction from spherical stereo image pairs. Computer Vision and Image Understanding (CVIU), 139(C):104–121, Oct. 2015, ISSN 1077-3142. http://dx.doi.org/10.1016/j.cviu.2015.04.001.  17, 69

# Results on Spherical Images

| RGB Image | Original Depth Map | Enhanced Depth Map | Input Volume | Predicted Volume |
|---|---|---|---|---|



floor █ wall █ window █ chair █ table █ sofa █ furn. █ objects █

## Conclusions

- We introduced the $360^o$ Semantic Scene Completion

- Works with high-end sensors or off-the-shelf $360^o$ cameras

- Segmentation accuracy equivalent to limited view solutions

- High levels of completion of occluded regions

# Publication

*Sematic Scene Completion from a Single 360° Image and Depth Map*



*\*Published in the proceedings of the 15th International Conference on Computer Vision Theory and Applications (VISAPP2020) (Qualis A1)*

[31] Dourado, A., Kim, H., de Campos, T.E., and Hilton, A.: Semantic scene completion from a single 360-degree image and depth map. In Proceedings of the 15th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISIGRAPP 2020), vol. 5: VISAPP, pp. 36–46. 7, 61

# Application Paper

*Immersive Audio-Visual Scene Reproduction using Semantic Scene Reconstruction from 360O Cameras*



https://www.cvssp.org/hkim/paper/CVST2020/

# Chapter 6

Work Plan

## Remaining Activities

- Review the most recent works on the subject
  - evaluate possible ways to improve EdgeNet (Chapter 4)
- Missing experiments:
  - try an offline very late fusion approach
  - train the $360^O$ solution on Stanford and other $360^O$ datasets (Chapter 5)
  - Try domain adaptation
    - from synthetic data
    - from NYUDV2
- Consolidate enhanced Chapters 4 and 5 into a Journal submission

# Timeline



| | 2018 | | 2019 | | 2020 | | 2021 | |
|---|---|---|---|---|---|---|---|---|
| | Jan-Jun | Jul-Dec | Jan-Jun | Jul-Dec | Jan-Jun | Jul-Dec | Jan-Jun | Jul-Dec |

Attend Classes
Surrey Internship
Skin Project
Edgenet Project
360° Project
Polishing Papers
*Attending VISAPP 2020 and UK Talk*
Qualif. Doc. Writing
Qualif. Doc. Internal Review
*Send Qualif. Doc. to Eval. Board*
NYUDv2 Video sequences experiments
Final Paper Writing and Submission
*Qualification Exam*
Thesis Writing
*Attending ICPR 2020*
Thesis Internal Review
*Send Thesis to Eval. Board*
*Thesis Presentation*

TODAY

# Thank you!

# Results – ablation study on SUNCG

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | **70.3** |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

# Results – ablation study on SUNCG

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | **70.3** |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

Effect of our efficient training pipeline

# Results – ablation study on SUNCG

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|-------|-------|------|------|-----|------|-------|------|------|-------|------|------|-------|------|------|------|------|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | **70.3** |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

Effect of our u-shaped architecture, with 3D dilated residial modules

# Results – ablation study on SUNCG

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | 70.3 |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

Effect of adding edges

# Results – ablation study on SUNCG

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | **70.3** |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

Effect of adding edges

# Results on NYU-DV2

| input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| d | SSCNet[24] | 76.3 | **95.2** | 73.5 | 96.3 | 84.9 | 56.8 | 28.2 | 21.3 | 56.0 | 52.7 | 33.7 | 10.9 | 44.3 | 25.4 | 46.4 |
| | SSCNet* | 92.7 | 89.7 | 83.8 | 97.0 | 94.6 | 74.3 | 51.1 | 43.7 | 78.2 | 70.9 | 49.5 | 45.2 | 61.0 | 51.3 | 65.2 |
| | DCRF [25] | – | – | – | 95.4 | 84.3 | 57.7 | 24.5 | 28.2 | 63.4 | 55.3 | 34.5 | 19.6 | 45.8 | 28.7 | 48.8 |
| | VVNetR-120 [9] | 90.8 | 91.7 | 84.0 | **98.4** | 87.0 | 61.0 | 54.8 | 49.3 | 83.0 | **75.5** | 55.1 | 43.5 | 68.8 | 57.7 | 66.7 |
| | EdgeNet-D | 93.1 | 90.4 | 84.8 | 97.2 | 94.4 | 78.4 | 56.1 | 50.4 | 80.5 | 73.8 | 54.5 | 49.8 | 69.5 | 59.2 | 69.5 |
| d+s | SNetFuse[14] | 56.7 | 91.7 | 53.9 | 65.5 | 60.7 | 50.3 | 56.4 | 26.1 | 47.3 | 43.7 | 30.6 | 37.2 | 44.9 | 30.0 | 44.8 |
| | TNetFuse[14] | 53.9 | 95.2 | 52.6 | 60.6 | 57.3 | 53.2 | 52.7 | 27.4 | 46.8 | 53.3 | 28.6 | 41.1 | 44.1 | 29.0 | 44.9 |
| d+e | SSCNet-E | 92.8 | 89.6 | 83.8 | 97.0 | 94.5 | 74.6 | 51.8 | 43.9 | 77.0 | 70.8 | 49.3 | 49.2 | 62.1 | 52.0 | 65.7 |
| | EdgeNet-EF(Ours) | **93.7** | 90.3 | **85.1** | 97.2 | 94.9 | **78.6** | 57.4 | 49.5 | 80.5 | 74.4 | **55.8** | 51.9 | 70.1 | **62.5** | **70.3** |
| | EdgeNet-MF(Ours) | 93.3 | 90.6 | **85.1** | 97.2 | **95.3** | 78.2 | **57.5** | **51.4** | **80.7** | 74.1 | 54.5 | **52.6** | **70.3** | 60.1 | 70.2 |
| | EdgeNet-LF(Ours) | 93.0 | 89.6 | 83.9 | 97.0 | 94.6 | 76.4 | 52.0 | 44.6 | 79.8 | 71.5 | 48.9 | 48.3 | 66.1 | 55.9 | 66.8 |

Effect of different fusion strategies

# Results on NYU-DV2

| train | input | model | scene completion | | | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | prec. | rec. | IoU | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| SUNCG | d | SSCNet[24] | 55.6 | 91.9 | 53.2 | 5.8 | 81.8 | 19.6 | 5.4 | 12.9 | 34.4 | 26 | 13.6 | 6.1 | 9.4 | 7.4 | 20.2 |
| | d+e | EdgeNet-EF(Ours) | **61.9** | 80.0 | **53.6** | 9.1 | **92.9** | 18.3 | 5.7 | 15.8 | 40.4 | 30.7 | 9.2 | 3.3 | 13.7 | 11.6 | 22.8 |
| | | EdgeNet-MF(Ours) | 60.7 | 80.3 | 52.8 | **11.0** | 92.3 | **20.5** | 7.2 | **16.3** | 42.8 | **32.8** | **10.5** | **6.0** | **15.7** | **11.8** | **24.3** |
| | | EdgeNet-LF(Ours) | 59.9 | **80.5** | 52.3 | 3.2 | 87.1 | 19.9 | **8.6** | 15.4 | **43.5** | 32.3 | 8.8 | 4.3 | 13.7 | 10.0 | 22.4 |
| NYU | d | SSCNet[24] | 57.0 | **94.5** | 55.1 | 15.1 | 94.7 | 24.4 | 0.0 | **12.6** | 32.1 | 35.0 | **13.0** | **7.8** | 27.1 | 10.1 | 24.7 |
| | d+e | EdgeNet-EF(Ours) | **78.1** | 65.1 | 55.1 | **21.8** | 95.0 | 27.3 | **8.4** | 6.8 | **53.1** | 38.6 | 7.5 | 0.0 | 30.4 | **13.3** | 27.5 |
| | | EdgeNet-MF(Ours) | 76.0 | 68.3 | **56.1** | 17.9 | 94.0 | **27.8** | 2.1 | 9.5 | 51.8 | **44.3** | 9.4 | 3.6 | **32.5** | 12.7 | **27.8** |
| | | EdgeNet-LF(Ours) | 75.5 | 67.5 | 55.4 | 19.8 | 94.9 | 24.4 | 5.7 | 7.2 | 50.3 | 38.8 | 10.0 | 0.0 | 33.2 | 12.2 | 27.0 |
| SUNCG + NYU | d | SSCNet[24] | 59.3 | 92.9 | 56.6 | 15.1 | 94.6 | 24.7 | 10.8 | 17.3 | 53.2 | 45.9 | 15.9 | 13.9 | 31.1 | 12.6 | 30.5 |
| | | DCRF[25] | - | - | - | 18.1 | 92.6 | 27.1 | 10.8 | 18.8 | 54.3 | 47.9 | 17.1 | 15.1 | 34.7 | 13.0 | 31.8 |
| | | VVNetR-120[9] | 69.8 | 83.1 | 61.1 | 19.3 | 94.8 | 28.0 | 12.2 | 19.6 | 57.0 | 50.5 | 17.6 | 11.9 | 35.6 | 15.3 | 32.9 |
| | d+c | Guedes et al. [7] | - | - | 56.6 | - | - | - | - | - | - | - | - | - | - | - | 30.5 |
| | d+s | Garbade et al. *[6] | 69.5 | 82.7 | **60.7** | 12.9 | 92.5 | 25.3 | 20.1 | 16.1 | 56.3 | 43.4 | 17.2 | 10.4 | 33.0 | 14.3 | 31.0 |
| | | SNetFuse[14] | 67.6 | **85.9** | **60.7** | 22.2 | 91.0 | 28.6 | **18.2** | 19.2 | 56.2 | 51.2 | 16.2 | 12.2 | 37.0 | 17.4 | 33.6 |
| | | TNetFuse[14] | 67.3 | 85.8 | **60.7** | 17.3 | 92.1 | 28.0 | 16.6 | 19.3 | **57.5** | **53.8** | **17.7** | **18.5** | **38.4** | **18.9** | **34.4** |
| | d+e | EdgeNet-EF(Ours) | 77.0 | 70.0 | 57.9 | 16.3 | **95.0** | 27.9 | 14.2 | 17.9 | 55.4 | 50.8 | 16.5 | 6.8 | 37.3 | 15.3 | 32.1 |
| | | EdgeNet-MF(Ours) | **79.1** | 66.6 | 56.7 | **22.4** | 95.0 | 29.7 | 15.5 | **20.9** | 54.1 | 53.0 | 15.6 | 14.9 | 35.0 | 14.8 | 33.7 |
| | | EdgeNet-LF(Ours) | 77.6 | 69.5 | 57.9 | 20.6 | 94.9 | 29.5 | 9.8 | 18.1 | 56.2 | 50.5 | 11.4 | 5.2 | 35.9 | 15.3 | 31.6 |

# Our approach

- Input volume:
  - 480 x 144 x 480 voxels
  - Voxel size: 0.02m
  - coverage: 9.6 x 2.8 x 9.6 m

- 8 partitions, emulating the field of view of a standard RGB-D sensor

- The partitions are taken from the sensor position, using a 45$^o$ step

- We move the point-of-view 1.7m back from the original sensor position, to get more overlapped coverage
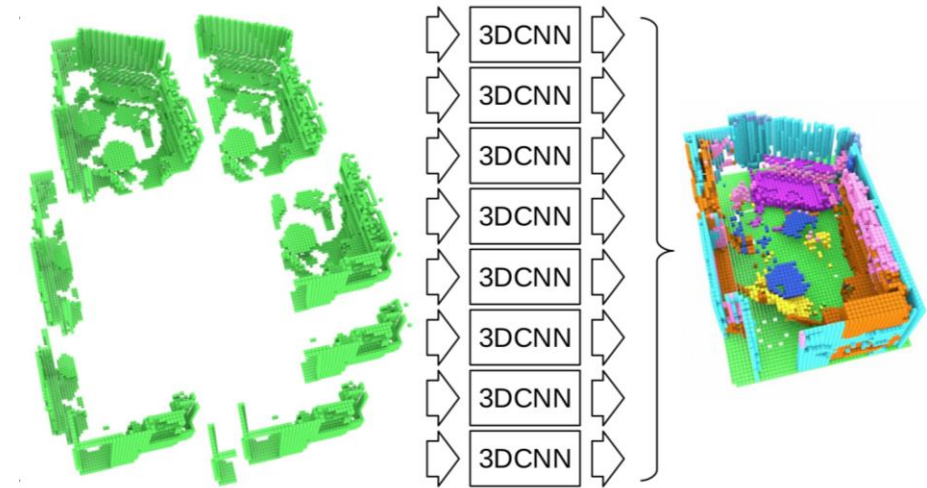


Input Partitioning

# Our approach

- Each partition of the input is processed by our CNN, generating 8 predicted volumes

- Overlapping areas are ensembled using the sum rule

- Each predicted partition size is 60 x 36 x 60

- The resulting ensembled volume size is 120 x 36 x 120



Prediction Ensemble

# Results on Stanford 2D-3DS Dataset

| evaluation dataset | model | scene coverage | semantic scene completion (IoU, in percentages) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | ceil. | floor | wall | win. | chair | bed | sofa | table | tvs | furn. | objs. | avg. |
| NYU v2 RGB-D | SSCNet | partial | 15.1 | 94.6 | 24.7 | 10.8 | 17.3 | 53.2 | 45.9 | 15.9 | **13.9** | 31.1 | 12.6 | 30.5 |
| | SGC | | 17.5 | 75.4 | 25.8 | **6.7** | 15.3 | 53.8 | 42.4 | 11.2 | 0.0 | 33.4 | 11.8 | 26.7 |
| | EdgeNet | | **23.6** | **95.0** | 28.6 | **12.6** | 13.1 | **57.7** | **51.1** | 16.4 | 9.6 | **37.5** | 13.4 | 32.6 |
| Stanford 2D-3D-S | **Ours** | full (360°) | 15.6 | 92.8 | **50.6** | 6.6 | **26.7** | - | 35.4 | **33.6** | - | 32.2 | **15.4** | **34.3** |