

Otto-von-Guericke University Magdeburg

Faculty of Computer Science



Master's Thesis

Synth2Real : 3D-Furniture Reconstruction in Ersatz Environment (S2R:3D-FREE)

Author:

Kartik Prabhu

October 15, 2021

Advisors:

Prof. X

National Aeronautics and Space Administration

Dr. No

Department of Technical and Business Information Systems

Otto-von-Guericke University Magdeburg

Prabhu, Kartik:

*Synth2Real : 3D-Furniture Reconstruction in Ersatz Environment
(S2R:3D-FREE)*

Master's Thesis, Otto-von-Guericke University Magdeburg, 2021.

Abstract

Besides the title, the abstract is the most important part of your thesis, as most readers will only read title and abstract. Your goal is to advertise the rest of the thesis for potential readers. For that, you briefly explain what you are focusing on in the thesis. With a misleading abstract, you will miss interested readers and maybe even attract readers with wrong assumptions about your work which will stop reading soon. The abstract should describe the general area as well as the most interesting insights of the work. It is crucial to find the right level of abstraction and length. An abstract typically consists of one paragraph that is significantly shorter than the introduction.

Abstracts typically follow the same structure. You start by describing the research area as well as the general and the specific problem you are focusing on. Then, you outline how you approach the problem in terms of concepts and evaluations. Finally, you close with the most interesting insights that you gained and why they are relevant for the research area.

Acknowledgments

If you want to, you can thank your advisors and/or anyone else who supported you during the thesis in some way.

Contents

List of Figures	ix
List of Tables	xi
List of Code Listings	xiii
List of Acronyms	xv
List of Symbols	xvii
1 Introduction	1
1.1 Domain Adaptation vs Domain randomisation	2
1.2 Volumetric representations of 3D shapes	3
1.3 GameEngines	3
1.4 Background and motivation	4
1.5 Goal of this Thesis	4
1.6 Structure of this Thesis	5
2 Related Work	7
2.1 Indoor datasets	7
2.1.0.1 Indoor synthetic datasets	7
2.1.0.2 Tools to create synthetic datasets	8
2.2 State of the art for 3d-reconstruction	8
2.3 Domain adaptation	9
2.3.1 GAN based style transfer	9
3 Concept	11
3.1 Pix3D: A large-scale benchmark	11
3.1.1 Disadvantages of Pix3D	11
3.1.2 Why Pix3D?	12
3.2 Role of SceneNet	12
3.3 Unity based pipeline	12
3.3.1 Domain Randomisation with Unity Engine	13
3.3.1.1 Scenes from scenenet	13
3.3.1.2 Camera Viewpoints	13
3.3.1.3 Lighting and shadows	14
3.3.1.4 Randomised textures	14
3.3.1.5 Replacing target objects	14

3.4	S2R:3D-FREE, a Pix3D based sythetic dataset	14
3.5	3D reconstruction pipeline	14
3.5.1	Pix2Vox	15
3.5.1.1	Why Pix2Vox?	15
3.6	Domain adaptation with Gan <i>if we decide to go with this approach</i>	15
4	Implementation	17
4.1	3D-Scene framework	17
4.2	3d-Reconstruction framework	17
4.3	Implement a domain adaptation technique	17
5	Evaluation	19
5.1	Domain gaps - Qualitative	19
5.2	Domain gaps - Quantitative	19
5.3	Performance	19
5.4	Domain shift technique	19
5.5	Ablation study on chairs	19
6	Conclusion	21
6.1	Summary	21
6.2	Limitations of game engines	21
6.3	Future Work	21
	Appendix	23
B	Tutorial	25
B.1	This is a section	25
B.1.1	This is a subsection	25

List of Figures

1.1	A graph showing how synthetic dataset will evolve overtime ?	2
1.2	3D representation of Stanford bunny model.(left to right) Point cloud, voel and mesh ?	3
3.1	Distribution of pix3D ? images(left), unique models(right)	12
3.2	Sample RGB images from pix3D	13
3.3	Network architecture for pix2vox ?	15
3.4	A survey conducted by ?, proves that Pix2Vox is considerably a good 3D reconstruction model. <i>replace this figure with better one</i>	16
A.1	Caption in the list of figures	23
B.1	List	26
B.2	List, Ordered	26
B.3	List, Set	26
B.4	List, Ordered, Set	26
B.5	Drawing with tables and TikZ	26

List of Tables

B.1 Short caption for list of tables 27

List of Code Listings

List of Acronyms

FOL	First-Order Logic
SLR	Systematic Literature Review
SPL	Software Product Line

List of Symbols

General

$\alpha[x \setminus y]$ Substitution of x in α with y

$\mathcal{P}(A)$ Power set of a set A

$f|_{A'}$ Restriction of a function $f: A \rightarrow B$ to a smaller domain $A' \subseteq A$

First-Order and Dynamic Logic

$A \wedge B$ Formalization of the sentence “ P and Q ”

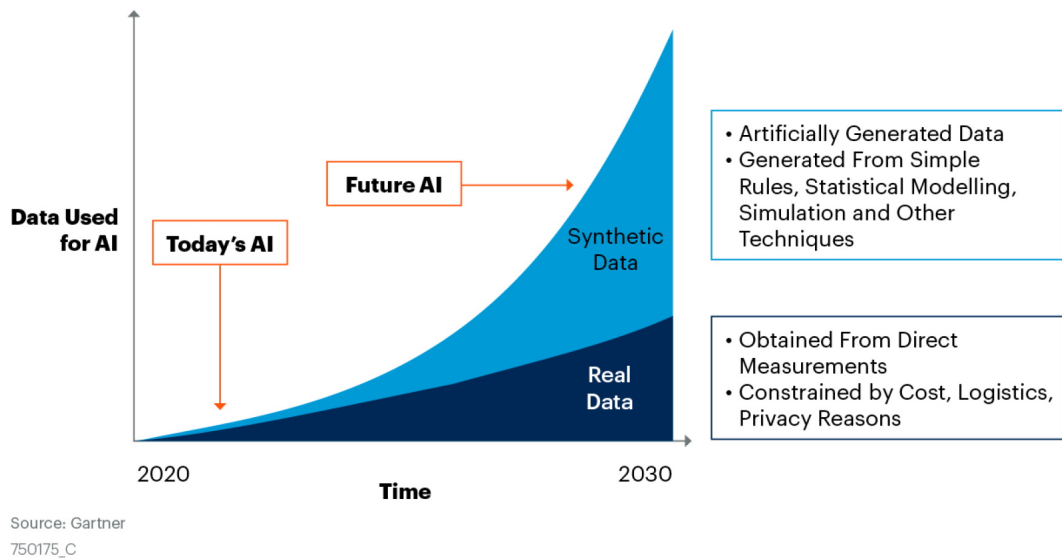
1. Introduction

As Deep Learning field goes deeper and wider, the need for abundant data has become the basic requirement. The performance of computer vision based on deep learning models heavily depends on the consistency of labelled images. But in practice, collecting data manually is time-consuming, expensive and needs lot of effort. As a solution, we can train models on synthetic dataset and deploy the models to real world problems. Unfortunately it is seen that with supervised learning, the model does not generalise when the distributions of dataset diverges and hence a model trained on synthetic dataset may not perform consistently with a real-world data.

Synthetic data can be defined as data created artificially using computer simulations or algorithms and not from real world events. These data can be created in abundance and hence are frontrunners for Deep Learning training. Synthetic approaches are also means of getting data which might be difficult to collect and annotate in real world. Further more, companies have realised that synthetic data are a necessity not just for task which are inconvenient to annotate, but also for improving the performance of existing models. In a research article published by Gartner [1], they claim that by 2030, synthetic datasets will overshadow real dataset in all domains. The progression of synthetic data generation is predicted to be as shown in 1.1

Deep Learning has made leaps and bounds in solving 2D-image tasks like classification, segmentations, object detection, etc. It has now entered the realm of Three-Dimension. An image is a 3D objects projects onto a 2D-surface. While doing so, some data from higher dimension is being lost in lower dimension. This inverse ability to reconstruct 3D objects from 2D images has a wide scale application in computer vision and robotics. Tasks like human body shapes and face reconstructions [2], 3D-Scene reconstructions [3], generic object mesh reconstruction with textures are some of the advancements we had over recent years. Here comes the challenge; 3D data is not easy to collect. 3D data is not just limited but also noisy and expensive to collect. The results of models depend on the consistency of the labelled data. In this thesis we discuss how game engines can contribute to generating quality labeled dataset for 3D reconstruction task.

By 2030, Synthetic Data Will Completely Overshadow Real Data in AI Models



Gartner

Figure 1.1: A graph showing how synthetic dataset will evolve overtime ?

1.1 Domain Adaptation vs Domain randomisation

As mentioned earlier, the success of Convolutional Neural Networks depends highly on the training data. As it so happens, not every task will have abundant data to train with. One solution for this problem is using a pre-trained network of similar problem and using transfer learning. The pre-trained network is then trained over with the limited data available for the problem at hand. In this case the features learnt from larger dataset is fine-tuned by overwriting few layers using small annotated dataset of a different domain. Domain adaptation is field of study in machine learning where the distribution of training data differs from testing or target distribution. One shallow approach is as mentioned above; re-weighting the network with testing data to adapt to the problem domain ?. *add more citations for domain adaptation.* "Deep domain adaptation methods leverage deep networks to learn more transferable representations by embedding domain adaptation in the pipeline of deep learning" ? For this thesis, domain adaptation is relevant as synthetic dataset might not represent the real or target domain, which can lead to drop in performance.

Domain randomisation on the other hand, is improving data quality such that the target domain distribution is included in the training domain. This is mainly used to bridge the 'reality gap' as termed in ?, between simulated environment and real images. This is achieved by randomising the rendering of objects in the synthetic data. The key concept here is that by rendering objects with large variance, the real-world distribution will appear to be just a variation in the dataset. Some of the features that can be randomised in a synthetic dataset are textures, brightness, shadows, camera settings, object pose, etc. This will be further discussed in chapter *mention chapter.*

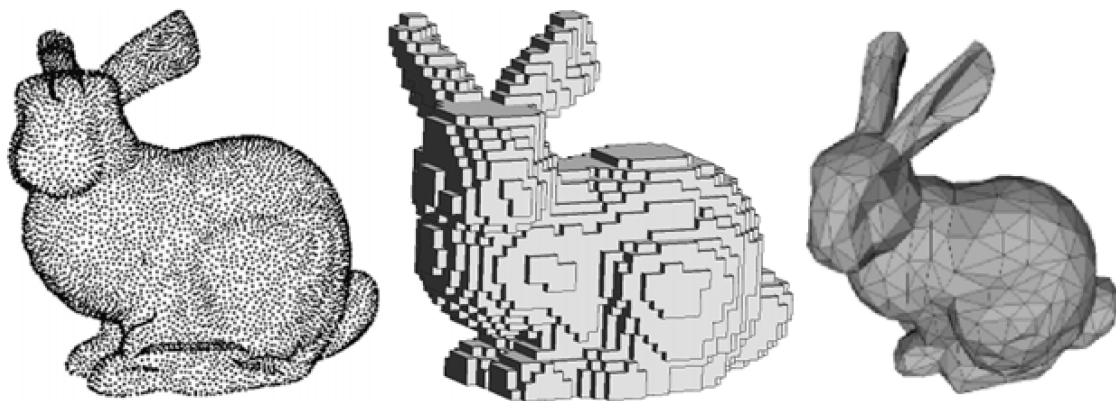


Figure 1.2: 3D representation of Stanford bunny model.(left to right) Point cloud, voxel and mesh ?

1.2 Volumetric representations of 3D shapes

The universal representation of images in computer graphics is via pixels. But unlike 2D, 3D data can be represented in various formats like voxels, mesh and point clouds. A sample of these 3 formats is shown in 1.2. Each of the representations has its own sets of advantages and disadvantages. Voxels, which is short for Volumetric Pixels is a 3D-grid of pixels of constant size. As indicated in ?, the main advantage of voxels is that Convolution Neural Networks can be easily applied to 3D as in 2D. But since most of the 3D geometry is surface based, it can be wasteful and computationally expensive. Mesh is composed of vertices, edges and faces in 3D space that indicate the formation of 3D objects. This form of representation is far more compact at granular level depending on the resolution. Point clouds on the other hand are collection of 3D points with (x,y,z) coordinates on the surface of the object. The collection of points determines detailing the 3D object representation. But in both these cases CNN is not directly applicable. In this thesis we will be dealing with voxel based models as the focus is on performance of synthetic data rather than the model or the 3d representation itself. *add an image with 3 types of representations.*

1.3 GameEngines

As per ?, game engine is defined as “A series of modules and interfaces that allows a development team to focus on product gameplay content, rather than technical content”. In layman terms game engines are software used for development of games. Unreal Engine, Unity, GameMaker, Amazon Lumberyard, CryEngine are some of the popular game engines in modern times. Over time, the graphics in each of the game engines has improved to an extent where visually we are not able to differentiate between reality and graphics. These use object-oriented techniques which help in modularity. Also, they are GUI-oriented and hence are user friendly and accessible even to a novice programmer. The key component in game engines which we will be focusing on is the 'Rendering engine'. 3D Rendering is process of converting 3D data into 2D images. Rendering can either be real-time or offline/pre-rendering. Game engines aim to make photorealism at highest level with least possible rendering time.

1.4 Background and motivation

3D reconstruction is of great significance to understand scene and to cross the 2D realm. But the task is not trivial, considering the input is still a Two-dimensional RGB image. It is even more challenging since obtaining new data is expensive and time-consuming. The 3D reconstruction field has datasets like ShapeNet [1], 3D-Front [2]. But there are no corresponding real-images to check if the model will work in real world scenarios. Also searching for real-world examples or creating a set-up scene and then labelling the images manually can lead to human mistakes resulting in wrong labels. A solution to this problem is automating the task of generating synthetic data which can produce pixel-perfect annotations. This can solve human-in-the-loop problem, and while time gets saved, large quantity of data can be generated. But quantity should not be the only focus while creating the data, equal importance needs to be given to quality of the data. Game engines have come long way from just supporting 2D to supporting games with photorealistic rendering. Even though game engines have grown to be successful, not many tools are based on it. As 3D models available, we can import these models into the game engine and create an ersatz environment to generate photorealistic dataset. Further, we can check if these photorealistic synthetic data can be used for 3D reconstruction task.

1.5 Goal of this Thesis

The following research questions will be answered:

1. Are game engines a good medium to create photorealistic synthetic dataset?
 - (a) *Experiment: Conduct a survey on real vs synthetic. Use dataset like Pix3d(real), blender proc, unity, 3d-front (synthetic)*
 - (b) *Depending on the survey, the most real dataset will be considered as baseline and other datasets will be compared to it.*
2. Can Ersatz environment from a game engine like Unity replace real data for training in 3d reconstruction task?
 - (a) *Train on s2r:3dfree data, test on pix3d*
3. To what extent can the performance of model Pre-trained with real dataset improved with Synthetic dataset from game engine, with size of synthetic dataset being 10(*depends on how much data is generated*) times that of real dataset?
 1. *pretrained on pix3d, continue training on s2r:3dfree*
 2. *Mix training*

(optional: Can the photorealism of synthetic image further be enhanced using domain adaptation/GAN techniques?)

Overall dataset experiments:

1. Ablation study on chairs, conduct all above mentioned experiments only on chairs as it has highest number of 3d models available
2. Finetune chair dataset and check if dataset can be improved to get better results
3. Scoremaps on voxel output to see if the probability of voxel is high in required region and to see which region is least to give possible explanation for low IOU

1.6 Structure of this Thesis

The rest of this thesis is structured as follows:

- *Chapter 2*, we visit some of the related work on synthetic data generation and state of the art 3d reconstruction networks.
- *Chapter 3* deals with concepts and design choices. x
- In *Chapter 4*, we discuss in detail the 3DScene tool developed using Unity game engine for creating ersatz environment.
- *Chapter 5* is dedicated to reviewing and discussing evaluation results obtained by comparison of synthetic and real dataset.
- Finally in *Chapter 6*, we conclude this thesis with results of the study, highlight few limitations and discuss future improvements.

2. Related Work

2.1 Indoor datasets

Indoor scene dataset has been in the rise with increasing interest in scene processing understanding ??????????. Synthetic dataset is not something new in the world of machine learning. As researchers realised the disadvantages of real dataset, focus was shifted to synthetic dataset. While ?? are real-world datasets, ??? are synthetically produced. The real-world dataset are gathered from live scans. The synthetic dataset can either be manually configured by a professional or automated by a programmer.

2.1.0.1 Indoor synthetic datasets

Alibaba group introduced 3D-FRONT ? which stands for 3D Furnished Rooms with layOuts and semaNTics dataset which comprises of synthetic indoor scenes designed under the supervision of professionals. It consists of 18,968 rooms and 13,151 textured furnitures. SceneNet ? is a large collection of photorealistic images and trajectories. This is discussed in detail in 3.2 section.

SunCG ? was a key dataset for scene understanding. The dataset contained over 45,000 variations of scenes with realistic room layout created manually. Each scene was semantically labeled and also provided with volumetric ground truth data. The dataset has also been used for task like depth estimation, semantic scene completion, SLAM, indoor navigation, etc. Unfortunately due to legal issues ¹ the dataset has been made publically unavailable which has left a void in the field.

Structure3D ? is another impressive synthetic data for indoor scenes which introduced their own photorealistic renderer. The dataset comprises 21,835 rooms in 3,500 scenes and 196k 2D-images rendered with photo-realism. But the CAD models of the 3D furnitures which is used to populate the scenes is not made available to public. And hence 3D reconstruction related tasks can not be performed. It is also demonstrated that with combination of synthetic and real dataset deep learning task for room layout estimation improved performance on benchmark datasets. This dataset is

¹<https://futurism.com/tech-suing-facebook-princeton-data>

more focused on room layout estimation and not 3D reconstruction, but with few modification, it can be mapped to 3d furniture reconstruction tasks.

Openrooms ? use *Scannet* ? as their layout foundation, retrieve corresponding models from *shapenet* ? and then replace the CAD model with retrieved model with proper alignment. They further add reflectance and illumination to compose photo-realistic images. As of August 2021, only the dataset has been made public and not the generation tool. Though the underlying concept of *Openrooms* is to convert existing scans into photo-realistic synthetic images, in terms of the output images we consider them our counterpart, as the framework can produce normals, depth maps, instances segmentation and masks same as we do.

HyperMism ? is Apple's repository for holistic indoor scene understanding. It is a collection of synthetic scenes created with the help of professional artist. ? was the starting point for the dataset for which assets were purchased from ?. The dataset includes images, 3D assets, semantic instance segmentations, and a disentangled image representation with diffused lighting and shading. Even though the 3d triangle meshes for each asset is available online, we have to purchase them to create custom dataset. They also admit that the costs to generate the dataset is expensive {approximately \$57K ?}.

InteriorNet ? claims to be a photo-realistic indoor scene simulator with realistic lighting and scenes which change over time. The image dataset includes rgb, depth and semantic segmentations. Along with images, they also provide synthesized realistic trajectories at video-frame rate with various motion patterns. The simulator also supports scenes from ? and ? along with their own database.

Another simulated framework for visual research is *House Of inteRactions (THOR)* introduced in *AI2-Thor* ?. This is again a Agent focused photo-realistic dataset with the key factor being actionable objects so that agents can interact with the objects or manipulate them. The underlying renderer for this framework is Unity game engine. *RoboThor* ? is built upon *AI2-Thor* which consists of real scenes and its corresponding synthetic equivalent. This helps in the study of behavior of agents in real world when trained on synthetic data.

Habitat: A Platform for Embodied AI Research ?, is a photorealistic 3D simulation which can be used for training virtual agents for tasks like navigation, question answering, intruction foloowing. The paper introduces *Habitat-Sim* which renders scenes from *Matterport3d* ?, *Gibson* ?, *Replica* ? and some other datasets. The focus of the simulator is providing the agent with sensor data and allowing additional sensors as plugins. At the foundation level *Habitat-sim* uses *Magnum graphics middleware library* ² which supports cross-platform on various hardware configuration.

2.1.0.2 Tools to create synthetic datasets

BlenderProc: Reducing the Reality Gap with Photorealistic Rendering ?

NVIDIA Deep learning Dataset Synthesizer (NDDS)

SynthDet: An end-to-end object detection pipeline using synthetic data A table to explain diffs between various methods(only if 3D-FREE has more advantage)

²<https://magnum.graphics/>

2.2 State of the art for 3d-reconstruction

Voxel based, mesh based Pix2vox, OctNet, Mesh rcnn, Pixel2Mesh, Occupancy networks, etc

2.3 Domain adaptation

Discuss other methods used to mitigate domain shift Example: distance learning, subspace matching

2.3.1 GAN based style transfer

3. Concept

In this section we discuss about dataset and design choice, and the rationale behind reducing the domain gap.

3.1 Pix3D: A large-scale benchmark

In [1], a large-scale benchmark for 2D-3D alignment was introduced. The raw images from web search engines were collected and labelled keypoints were used to align the 2D images with the corresponding 3D shapes. The 3D models are extension from IKEA dataset [2], which is a collection of high-quality IKEA furniture. The dataset also provides with masks and keypoints for the object under observation. For adding more images to the IKEA dataset, the authors of [1] conducted manual web search on Google, Bing, and Baidu, using Ikea model name as keywords. to get around 104,220 images which were further filtered by removing irrelevant images with the help of Amazon Mechanical Turk (AMT) workers. After this manual experiment the total images in Pix3D, only 14,600 images were selected for the 219 Ikea models. For our experiment on synthetic to real dataset, we chose to select only furniture classes from Pix3D, leaving out "misc and tools" classes, which were significantly less to begin with.

3.1.1 Disadvantages of Pix3D

The distribution of models in pix3d is as shown in 3.1. As we can see, the dataset distribution is uneven across classes and more than 50% of classes have less than 1000 images.

Though Pix3D set a benchmark for 2D-3D alignment, here are few disadvantages of using this real dataset.

- 1. For Deep Learning approaches, we need large-scale data and 14,600 might not be sufficient.*
- 2. The orientation of object is not randomised.*
- 3. The dataset does not provide 2.5D information (i.e. depth and normal) which can be crucial for 3D learning.*

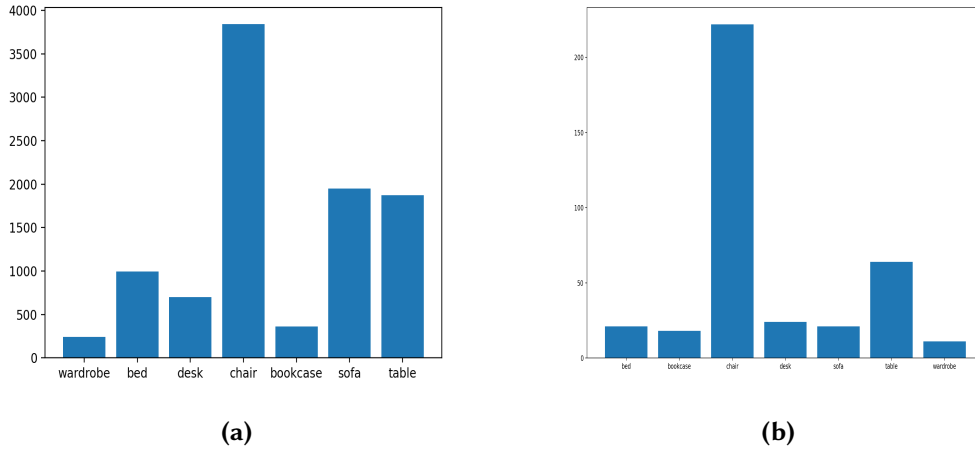


Figure 3.1: Distribution of pix3D ? images(left), unique models(right)

3.1.2 Why Pix3D?

Pix3D is a collection of indoor scenes with complex background, varying light conditions with shadows, reflective surfaces and even varying occlusion levels. Each image comprises of collection of objects in the scene, but only one object from the category is annotated. It is a perfect example of having limited real world data. And since 3D models are available for each furniture, a synthetic data can be generated in abundance from those models.

3.2 Role of SceneNet

SceneNet ? is a large collection photorealistic indoor scene trajectories. The dataset provides with images and videos of indoor scene which can be used for tasks like SLAM, semantic and instance segmentation, object detection and further enhanced for other vision problems like optical flow depth and pose estimation ?). They use ShapeNet ? models to occupy 57 indoor scenes which gives the scene unlimited configurations. Unfortunately there is no mapping from scene to 3D model for tasks like 3D reconstructions. In our approach, we utilize the scene provided by SceneNet as layout for our indoor scenes. This includes initial shapenet models with furniture placement. We then replace the class under observation with a corresponding model from Pix3D.

3.3 Unity based pipeline

Unity game engine is used for creating the synthetic dataset. Unity is a cross-platform game engine, the basic version of which is available for free. There is no official comparison between Unreal Engine and Unity engine to have a deciding factor. The selection of unity engine was purely individual choice and ease of use. But both these game engines have ability to create realistic looking scenes which can be used for synthetic dataset generation. With Unity game engine, the available scenes from scenenet and 3D models from pix3d can be imported to form an ersatz environment.



Figure 3.2: Sample RGB images from pix3D

Further domain randomisation can be applied to create a dataset of photorealistic images and 2.5D data like normal, depthmaps, masks.

3.3.1 Domain Randomisation with Unity Engine

In this section we discuss what kind of domain randomisation is applied for the dataset creation.

3.3.1.1 Scenes from scenenet

As discussed in section 3.2, SceneNet provides wide variety of indoor scenes which includes bedrooms, bathroom, living room, kitchen and office. For our case, we did not consider bathroom as the class under observations are furnitures and we rarely have any furnitures in the bathroom. For the rest of the scenes we have a pre-determined setup given by scenenet. We use 25 scenes in total as basic rooms for our target object to be place. And with further randomisation we feel this should be sufficient to check if unity can produce a useful dataset.

3.3.1.2 Camera Viewpoints

The position and orientation of camera gives us the most randomisation in this setup. The minimum and maximum distance can be chosen by the user for the camera to be places from the target object. A random point can be selected in this range as the position for the camera. As for the orientation of the camera, we can make the camera to always look at the target object, irrespective of the position. With this, we

achieved different background for target object, and the target object itself appears to be in different position and orientation.

3.3.1.3 Lighting and shadows

Lighting plays an important role in photorealism. If the lighting is not setup correctly the synthetic dataset will never seem to be photorealistic. Unity offers wide variety of lighting like global light which acts like the sunlight, and various indoor lighting systems. Ideally we should make the luminous objects like lamps, chandelier, bulbs etc be the source of light for indoor scenes, but we observed that the room does not lit up uniformly which makes it less photorealistic. Hence we use some pre-determined lighting settings which will be further discussed in implementation section. [mention section](#).

3.3.1.4 Randomised textures

SceneNet ? also provides with textures for different categories in the scene. We then further increase the texture database by adding few more textures from ambientCG.com, which provide license for free. We randomly allocate textures to each object making sure that each category has same texture to make the scene more uniform.

3.3.1.5 Replacing target objects

To further randomise the scene, the category of target objects in the scene are replaced by the object under observation. When there are more than one object of same category we randomise the selection of which object is to be replaced which further randomises the captured data. The scale of the target object is made to match the scale of the category object in the scene in the least dimension. For example if the length of the category object in the scene is the least among length, width and height, then the target object is scaled to match this length. This makes the target object blend in with the scene.

3.4 S2R:3D-FREE, a Pix3D based sythetic dataset

S2R:3D-FREE dataset, which stands for Synth2Real: 3-Dimensional Furniture Reconstruction from Ersatz Environment, is a combination of SceneNet and Pix3D dataset. We utilize the availability of 3D models of rooms and furnitures from these 2 dataset to create an ersatz environment using Unity as our framework. We randomise the indoor scenes from SceneNet ? along with textures provided by them with some additional complex textures. The other option would have been to use the shapenet as the target model, but in this case we would not be having a real dataset to compare to. Ideally a model trained on shapenet should also perform well with a real dataset like pix3d, but we decided to have a synthetic dataset based on pix3d itself for better comparison.

3.5 3D reconstruction pipeline

Now that we discussed on how the dataset will be created, we move to the deep learning aspect of this thesis.

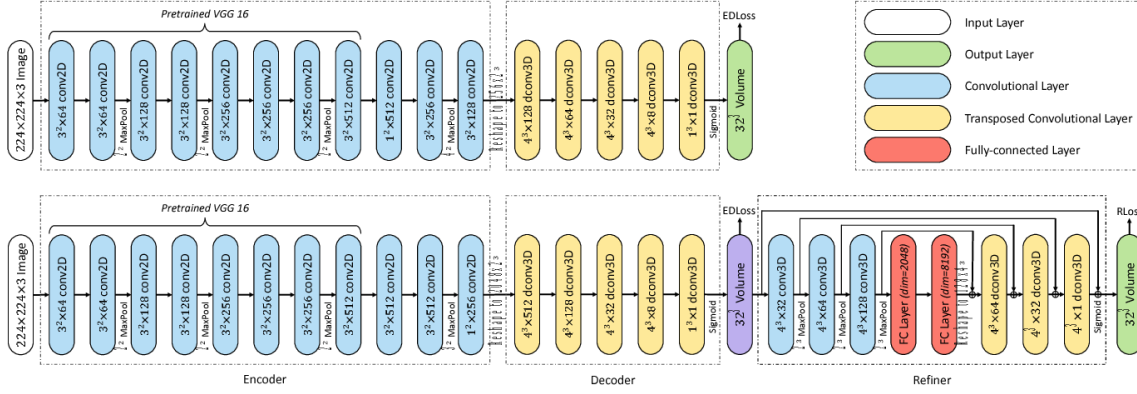


Figure 3.3: Network architecture for pix2vox ?

3.5.1 Pix2Vox

The architecture of pix2vox ? as in 3.3. The network is comprised of 4 modules: Encoder, Decoder, Merger and Refiner. Merger plays significant part when it comes to mult-view reconstruction of 3d objects. But we focus on single-view 3d reconstruction and hence merger will not influence the output significantly. As Encoder, pix2vox has utilised VGG16 network ? pre-trained on ImageNet ?. Decoder is an expansive network which converts 2D embeddings to 3D voxel grid. Refiner is an auto-encoder which takes the 3D output from the merger and produces more refined final output.

3.5.1.1 Why Pix2Vox?

Pix2Vox has been used as baseline by most of the research oriented to 3D reconstruction. This network is one of the few networks to be tested on pix3d dataset. According to the survey conducted by ?, performance of pix2vox ? has be noted to be significantly higher when compared to previous work *cite other papers* as shown in 3.4. While this comparison was made on 3d reconstruction of shapenet dataset, since pix3d was not available at the time when previous work were published. From our survey, only CoReNet ? had a slight gain in performance compared to pix2vox. When trained on shapenet and tested with pix3d, CoReNet gave a result of 29.7% IoU while Pix2Vox gave a result of 28.8% IoU and Pix2Vox++ a result of 29.2% IoU. Since the difference in the performance was not statistically significant we decided to stick with the baseline model itself. Another reason for selecting pix2vox model is that the backbone of the architecture is pre-trained with ImageNet and hence the embeddings generated from this encoder can help in visualising the domain space of both Pix3D (real images) and S2R:3DFREE(synthetic images). *write about pix2vox++, which uses resnet as backbone if used.* Furthermore, the focus of this thesis is not to check which is the best model to reconstruct the model, but to check if game engines can produce photorealistic images usable for 3d reconstruction. Hence the selection of the model was not of utmost importance.

3.6 Domain adaptation with Gan *if we decide to go with this approach*

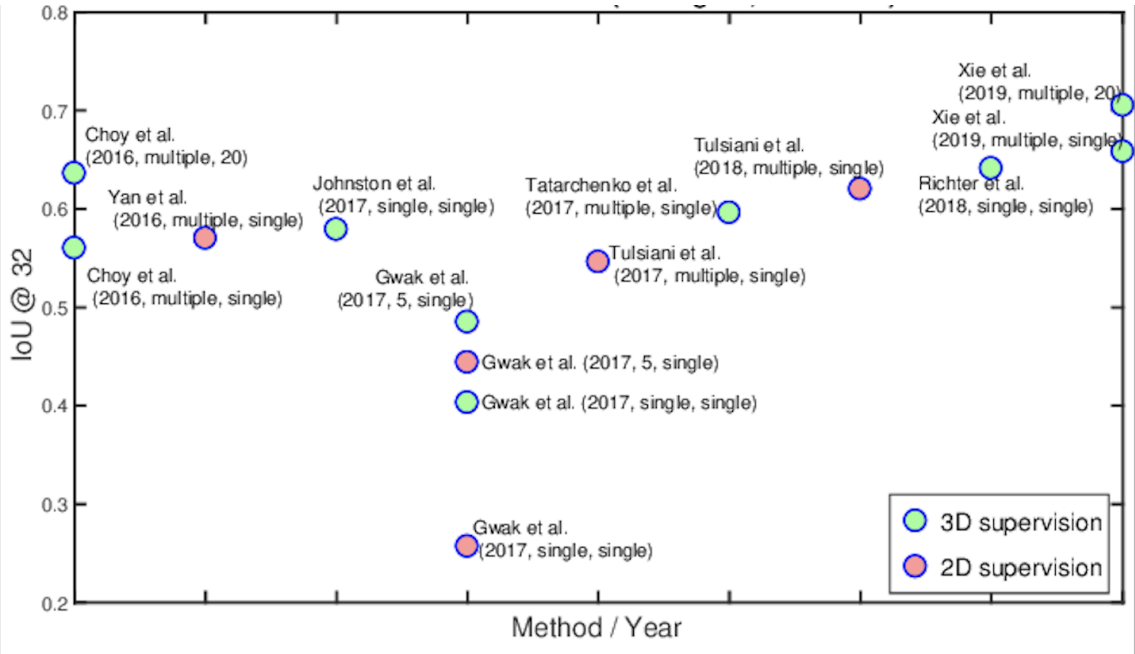


Figure 3.4: A survey conducted by ?, proves that Pix2Vox is considerably a good 3D reconstruction model. *replace this figure with better one*

4. Implementation

4.1 3D-Scene framework

Replacing models, lighting(indoor and outdoor), randomisation of texture, camera distance, camera angles, outdoor skybox, UI controls Manual inputs, automated Unity - ML-ImageSynthesis UML diagram is needed? (If it helps)

4.2 3d-Reconstruction framework

Development environment Training setup(hardware, voxel size...) Mixed precision using apex If mixed precision helped performance, memory, speed, evaluation metric

4.3 Implement a domain adaptation technique

To serve as a proof if more domain gap reduction is needed

5. Evaluation

5.1 Domain gaps - Qualitative

Visualize domain gap between S2R:3D-FREE and real image(pix3d) Visualize domain gap between other synthetic dataset like 3dfront, scenenet tsne

5.2 Domain gaps - Quantitative

Maximum Mean Discrepancy Kullback-Leibler divergence

5.3 Performance

*Performance of model trained on synthetic dataset by testing model with real dataset(pix3d)
Qualitative: Voxel output comparisons Quantitative: IOU, Dice*

5.4 Domain shift technique

Compare the differences in domain shift of models learnt on S2R:3D-FREE and a traditional domain shift learning.

5.5 Ablation study on chairs

*Check if randomising textures of chair and background helps improve performance.
Example: Create dataset with constant room texture vs randomising texture, Check if lighting helps (indoor, outdoor). Example: Create a dataset with constant lighting(only outdoor lighting) vs indoor lighting. Iou per category: as in <https://arxiv.org/pdf/1905.03678.pdf>*

6. Conclusion

6.1 Summary

6.2 Limitations of game engines

6.3 Future Work

Appendix

This table is automatically generated from a CSV file. Of course, you can also create tables from scratch.

	List_ Insert	List_ Search	Ordered_ Insert	Ordered_ Search	Ordered_ Min	Ordered_ Sort	Set_ Insert
coarse, Ex., Strict	2400	1418	10608	10824	2871	10169	23366
coarse, Ex., Def.	2400	1418	9652	20172	2871	14043	35884
coarse, Ex., None	2400	1418	14347	20160	2871	12687	54907
coarse, L.S., Strict	2400	1418	10608	10824	2871	10169	21047
coarse, L.S., Default	2400	1418	9430	20172	2871	13835	24688
coarse, L.S., None	2400	1418	13802	20160	2871	12687	28501
coarse, Complete	2400	1418	4504	10149	2871	10047	16114
coarse, Product	9600	2836	9010	20298	5742	20094	16114
fine, Ex., Strict	2211	1321	11210	16984	2501	6500	12901
fine, Ex., Def.	2211	1321	10299	16971	2501	6420	14695
fine, Ex., None	2211	1321	10881	16836	2501	7115	23005
fine, L.S., Strict	2211	1321	11210	16984	2501	6500	12301
fine, L.S., Default	2211	1321	10251	16971	2501	6420	13173
fine, L.S., None	2211	1321	11116	16836	2501	7115	20053
fine, Complete	2211	1321	2018	16531	2501	6613	9336
fine, Product	8844	2642	4036	33062	5002	13226	9336

Figure A.1: Caption in the text.

B. Tutorial

Acronyms like *software product lines (SPLs)* and *first-order logic (FOL)* are automatically added to the list of acronyms above. To reference a chapter, subsection, figure, etc., put a label after the command (see above) and then use *Chapter 1*.

You can mark stuff as **TODO**, which is useful for writing drafts.³

*Or put notes
to your
advisors in
the margin.*

Theorem B.1: My Cool Theorem

You can also add definitions and theorems in mathematical theses.

Proof. You can also add proofs, though most theses do not require any. □

In *Figure B.5*, we show an example of how to typeset a complex figure only with *LaTeX* by using *TikZ*.

B.1 This is a section

For citing, put entries in the *BibTeX* file (one example is there) and use the following to reference the authors directly: *? conducted a systematic literature review (SLR)* (this is an acronym defined in the main file), or use the following for references in parenthesis *[?]*. *? ?*

B.1.1 This is a subsection

Numerated list

1. *One*
2. *Two*
3. *Three*

³Footnotes are also possible, but are rarely used in computer science.

0.29

	<i>List</i>	<i>Ordered</i>	<i>Set</i>
Insert	●	○	○
Search	●	○	×
Sort	×	○	×

Figure B.1: List

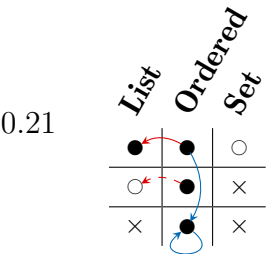


Figure B.2: List, Ordered

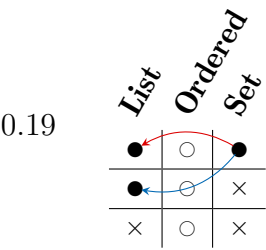


Figure B.3: List, Set

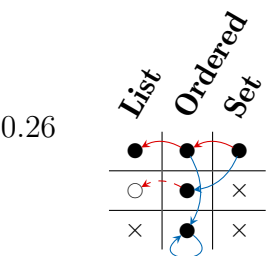


Figure B.4: List, Ordered, Set

Figure B.5: Drawing with tables and TikZ.

Table B.1: Long caption for a table.

Column1	Column2	Column3	Column4
Left	Right	Centered	Left, fixed width
Multi row		Multi column	
		X	
		Y	

Bullet list

- *One*
- *Two*
- *Three*

I herewith assure that I wrote the present thesis independently, that the thesis has not been partially or fully submitted as graded academic work and that I have used no other means than the ones indicated. I have indicated all parts of the work in which sources are used according to their wording or to their meaning.

Magdeburg, 15th October 2021