# Declaration

I hereby declare that the work presented in this thesis titled *Fine Grained Feature Representation using Computer Vision Techniques for Understanding Indoor Space* submitted to the Indian Institute of Technology Jodhpur in partial fulfillment of the requirements for the award of the degree of Doctor of Philosophy, is a bonafide record of the research work carried out under the supervision of Dr. Chiranjoy Chattopadhyay and Dr. Gaurav Bhatnagar. The contents of this thesis in full or in parts, have not been submitted to, and will not be submitted by me to, any other Institute or University in India or abroad for the award of any degree or diploma.

*Shreya Goyal*
*P16CS003*

# Certificate

This is to certify that the thesis titled *Fine Grained Feature Representation using Computer Vision Techniques for Understanding Indoor Space*, submitted by *Shreya Goyal (P16CS003)* to the Indian Institute of Technology Jodhpur for the award of the degree of *Doctor of Philosophy*, is a bonafide record of the research work done by her under our supervision. To the best of our knowledge, the contents of this thesis, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

*Chiranjoy Chattopadhyay*
*Ph.D. Thesis Supervisor*

*Gaurav Bhatnagar*
*Ph.D. Thesis Supervisor*

# Acknowledgments

I wholeheartedly thank my Ph.D. Thesis Supervisors, *Dr. Chiranjoy Chattopadhyay* and *Dr. Gaurav Bhatnagar*, for being mentor, guide and friend in the truest form during the entire course of my PhD programme. Qualities like patience, optimism and seeking perfection in one's work are just some of the wonderful things that I have learnt from them and would carry with me throughout my life. This dissertation would not have been possible without their constant support and guidance.

I thank Head of the Computer Science and Engineering Department, *Professor Richa Singh* for her guidance. I would also like to thank the members of the Doctoral Committee, *Dr. Anil Kumar Tiwari*, *Dr. Mayurakshi Chaudhury*, *Dr. Suman Kundu*, for their enthusiastic and continued guidance during the research work. I am also thankful to all the CSE department, IIT Jodhpur staff for their continuous and valuable support.

A big thanks to *Dr. Naimul Khan* from multimedia laboratory, Ryerson University, Toronto for providing me a welcoming and conducive working environment and all the brainstorming ideas and discussions during my visit to his laboratory.

My stay at the Institute was a wonderful experience because of my friends, *Abhra Paul*, *Devika Laishram*, *Rohit Shandley*, *Vishwa Deepak*, *Gaurav Bahuguna*, *Ajay Urgunde*. I thank them all for all the liveliness they infused into the non-academic part of the days at IIT Jodhpur.

I would like to thank my PhD colleagues and lab-mates *Divya Shrivastava*, *Arka Ujjal Dey*, *Dipti Trivedi* and *Jaspreet Kaur* for being the major source of support. A special mention to my friends *Shubhra Sharma* and *Anirudhha Singhal* for always being there for me with the words of encouragement.

Lastly, I deeply acknowledge and thank my parents and my brother *Aditya*, for their unconditional trust, endless patience, and giving me all the opportunities that have helped me shape into the individual that I am today.

# List of Figures

# List of Tables

# List of Symbols

| | |
|---|---|
| $I$ | Input floor plan image |
| $CC()$ | Connected Component Operator |
| $\zeta$ | Individual connected component |
| $\lvert.\rvert$ | Cardinality operator |
| $P_v$ | Pivot point |
| $\mathbb{W}$ | Formed word |
| $\mathbb{V}$ | Predefined word |
| $ED()$ | Euclidian distance operator |
| e | Edit distance |
| $W$ | Wall image |
| $\mathbb{E}$ | Edge image of walls |
| $W_H$ | Horizontal line image |
| $W_V$ | Vertical line image |
| $\mathbb{C}$ | Regions of intersection |
| $\mathbb{H}_B$ | Normalized histogram for Brick wall material |
| $\mathbb{H}_C$ | Normalized histogram for Concrete wall material |
| $\mathbb{H}_W$ | Normalized histogram for Wood wall material |
| $\mathbb{H}_S$ | Normalized histogram for given segments of wall materials |
| $\mathbb{S}$ | Set of wall segments |
| $\mathbb{D}_k$ | Distance between a pair of histograms of material |
| $(\beta, \xi)$ | Bins corresponding to histograms of current Wall segment |
| $\mathbb{F}_1$ | Wall material characterized floor plan image |
| $(\mathrm{Cx}_s, \mathrm{Cy}_s)$ | Centroid of each wall segment |
| $\mathrm{D}_p$ | Entry door pixel |
| $C$ | Connected Components |
| $\mathscr{A}_k$ | Area of each component |
| $A$ | Sorted collected of areas for all components |
| $F$ | UDI Signature |
| $\mathscr{R}_c$ | Room center |
| $\mathscr{D}_c$ | Decor Center |
| $d_n$ | Normalized distances for decor items |
| $D$ | Decor Items |
| $N_r$ | Number of rooms |
| $Dim_F$ | Dimension of the Feature Vector |
| $t$ | Shrinking Factor |
| $AM_D$ | Door based Adjacency Matrix |
| $AM_N$ | Navigation based Adjacency Matrix |
| $D_E$ | Door Entry |
| $D_X$ | Door Exit |
| $V_L$ | Vertex list |
| $C_S$ | Strongest Corners |
| $P^i$ | Navigation Route |
| $S_i$ | Sentences |
| $NN_r$ | neighboring Rooms |

| | |
|---|---|
| $DLOC$ | Relative location of rooms |
| $DIR$ | Direction |
| $N_m$ | number of turns |
| $R_i$ | Region Proposals |
| PV | Pooled Vector |
| $M$ | Projection Matrix |
| $Th$ | Threshold |
| $K$ | Paragraph Description |
| $Sent_{max}$ | Maximum Sentences |
| $Word_{max}$ | Maximum Words |
| $Prob$ | Probability term |
| $W_e$ | Word cues extracted from images |
| $C_i$ | Captions |
| $S_e$ | Target sentence embeddings |
| $h_t$ | Target hidden state |
| $h_s$ | Source hidden state |
| $e_{ij}$ | Alighnment scores |
| $\alpha_{ij}$ | Attention weights |
| $cv_i$ | Context vector |
| $r_i$ | Characterized Rooms |
| $o_i$ | Detected Objects |
| $I_k$ | Image for input indoor scene |
| $C^k$ | Camera pose |
| $f$ | focal length |
| $(C_x, C_y)$ | Camera center coordinates |
| $(\mathbf{u}, \mathbf{v})$ | Coordinates in in RGB image |
| $D_{u,v}$ | Depth value in the depth map |
| S | Scaling factor of the scene |
| $\mathbf{R}$ | Total number of rooms in dataset |
| $R_j$ | Individual rooms |
| $P_i$ | Point clouds |
| K | Boundary coordinates |
| $c_k$ | Clusters |
| $m_k$ | Mean for each cluster |
| $RP_i$ | Regularized local point cloud |
| $TP_i$ | Transformed point cloud |
| $\theta_x, \theta_y, \theta_z$ | Rotation angles |
| $(t_x, t_y)$ | Translation Coordinates |
| $FP$ | Final Polygon |
| p | Number of sides |
| $\phi$ | Angle |
| $line_i$ | Line joining two points |
| $s_i$ | Sides of polygon |
| $(X, Y, Z)$ | Coordinate system |
| $\mathbf{C}_{hull}$ | Convex hull |
| $\mathbf{Ratio}_D$ | Ratio used for door marking |
| $C_{BB_I}$ | Centroid of bounding box for door in real world image |
| $W_I$ | Wall in real world image |
| $L_{W_I}$ | Width of wall in real world image |
| $L_{W_F}$ | Width of wall in 2D mapping |
| $C_{BB_F}$ | Centroid of a 2D door symbol |

| | |
|---|---|
| $W_{I_F}$ | Wall in corresponding 2D mapping |
| $(\mu_x, \mu_y)$ | Intensity terms |
| $MAX_I$ | Maximum pixel value in image |
| $(\sigma_x, \sigma_y)$ | Standard deviation terms |

# List of Abbreviations

| | |
|---|---|
| **DIA** | Digital Image Analysis |
| **OCR** | Optical Character Recognition |
| **DAR** | Document Analysis and Research |
| **NLTK** | Natural Language Tookit |
| **LDA** | latent Dirichlet allocation |
| **JSON** | JavaScript Object Notation |
| **LBP** | Linear Binary Pattern |
| **NLP** | Natural Language Processing |
| **DIA** | Digital Image Analysis |
| **SLIC** | Simple Linear Iterative Clustering |
| **XML** | Extensible Markup Language |
| **ROUGE** | Recall-Oriented Understudy for Gisting Evaluation |
| **METEOR** | Metric for Evaluation of Translation with Explicit ORdering |
| **BLEU** | Bilingual Evaluation Understudy |
| **mAP** | Mean Average Precision |
| **IoU** | Intersection over Union |
| **BoD** | Bag of Decors |
| **LOFD** | Local Orientation and Frequency Descriptor |
| **BIM** | Building Information Modelling |
| **UDI** | Unique Decor Identifier |
| **BP** | Brevity Penalty |
| **SVM** | Support Vector Machine |
| **CNN** | Convolutional Neural Network |
| **TBDG** | Transformer Based Description Generation |
| **DSIC** | Description Synthesis using Image Cues |
| **RNN** | Recurrent Neural Network |
| **LSTM** | Long Short-Term Memory |
| **GRU** | Gated Recurrent Units |
| **RPN** | Region Proposal Network |
| **YOLO** | You Only Look Once |
| **F-RCNN** | Faster- Region based Convolutional Neural Network |
| **SSD** | Single Shot Detector |
| **FOV** | Field Of View |
| **HRNN** | Hierarchical Recurrent Neural Network |
| **SLAM** | Simultaneous Localization And Mapping |
| **AR** | Augmented Reality |
| **VR** | Virtual Reality |
| **ReLU** | Rectified Linear Unit |
| **CAD** | Computer Aided Design |
| **IMU** | Inertial Measurement Unit |
| **ROBIN** | Repository Of BuildIng plaNs |
| **SESYD** | Systems Evaluation SYnthetic Documents |
| **CVC-FP** | Computer Vision Center- FLoor Plan |

| | |
|---|---|
| **BRIDGE** | Building plan Repository for Image Description Generation, Evaluation |
| **SIFT** | Scale Invariant Feature Transform |
| **YOLO** | You Only Look Once |
| **IoU** | Intersection Over Union |
| **RCNN** | Region-based Convolutional Neural Networks |
| **SSD** | Single Shot Detector |
| **GUI** | Graphical User Interface |
| **PSNR** | Peak Signal to Noise Ratio |
| **SS** | Structural Similarity |
| **MSE** | Mean Square Error |
| **GT** | Ground Truth |
| **CE** | Corner Error |
| **PE** | Pixel Error |