# Contents

x

# List of Figures

# List of Acronyms

| | |
|---|---|
| **Adam** | Adaptive Moment Estimation (an optimisation algorithm) |
| **AI** | Artificial Intelligence |
| **Aleph** | A Learning Engine for Proposing Hypotheses (an ILP system) |
| **BotGNN** | Bottom-Graph Neural Network |
| **BK** | Background Knowledge |
| **CNN** | Convolutional Neural Network |
| **CONV** | Convolution (used for a block or a layer) |
| **DL** | Deep Learning |
| **DNN** | Deep Neural Network |
| **DRM** | Deep Relational Machine |
| **GNN** | Graph Neural Network |
| **ILP** | Inductive Logic Programming |
| **MDIE** | Mode-Directed Inverse Entailment |
| **ML** | Machine Learning |
| **MLP** | Multilayer Perceptron |
| **NCI** | National Cancer Institute |
| **NN** | Neural Network |
| **POOL** | Pooling (used for a block or a layer) |
| **RNN** | Recurrent Neural Network |
| **SGD** | Stochastic Gradient Descent |
| **SMILES** | Simplified Molecular-Input Line-Entry System |
| **VAE** | Variational Autoencoder |
| **VEGNN** | Vertex-Enriched Graph Neural Network |
| **XAI** | Explainable Artificial Intelligence |