

Applications of Deep Learning: Vision-Based Image Understanding and Reinforcement Learning in Gaming

Course: Spring 2025 – Advanced Machine Learning (BA-64061-002)

1. Summary

This report investigates two prominent applications of deep learning:

1. A vision pipeline using Conditional Generative Adversarial Networks (cGANs), Faster R-CNN for object detection, and BLIP for image captioning.
2. A reinforcement learning framework, theoretically designed, that uses a Deep Q-Network (DQN) to play the Flappy Bird game.

Part A was fully implemented and achieved a mean average precision (mAP) of 0.76 on object detection tasks. Part B is a conceptual design outlining the architecture, reward strategy, and expected learning dynamics of a DQN-based agent. Together, these parts demonstrate deep learning's strength in both perception-based and control-based tasks.

2. Introduction

Deep learning continues to drive significant advancements in artificial intelligence. It is central to enabling machines to perceive, generate, and act in complex environments. This report presents two contrasting but complementary use cases of deep learning:

- **Part A** focuses on visual understanding by combining a generative model (cGAN), a detection model (Faster R-CNN), and a captioning model (BLIP). This layered vision system simulates end-to-end visual intelligence, from content generation to semantic interpretation.
- **Part B** outlines a theoretical design of a reinforcement learning agent using DQN to play the Flappy Bird game. It demonstrates how an agent can learn to navigate and make decisions in a sequential, feedback-driven environment.

These projects represent real and conceptual implementations of deep learning, emphasizing its versatility across problem domains.

3. Current Research

Generative Adversarial Networks (GANs), proposed by Goodfellow et al. (2014), are a class of unsupervised learning models where two neural networks—the generator and discriminator—compete to produce realistic

data.

Conditional GANs (cGANs) improve upon this framework by including class labels to control the data generation process (Mirza & Osindero, 2014).

Faster R-CNN, developed by Ren et al. (2015), is a state-of-the-art object detection architecture that merges a Region Proposal Network (RPN) with convolutional neural networks to efficiently locate and classify objects in an image.

BLIP (Bootstrapping Language-Image Pretraining), introduced by Li et al. (2022), combines vision transformers with textual generation to produce contextual captions for detected visual elements.

Deep Q-Networks (DQN), introduced by Mnih et al. (2015), allow reinforcement learning agents to approximate Q-values using deep neural networks. They incorporate experience replay and target networks to stabilize training and are widely used in control-oriented tasks such as gameplay and robotics.

4. Data Collection and Model Development

Part A: Vision Pipeline Implementation

- **Dataset:** MS COCO (20,000 labeled images with bounding boxes and captions)
- **Image Generation:** A cGAN was trained to generate synthetic images conditioned on object categories from COCO
- **Object Detection:** A Faster R-CNN model pretrained on Open Images was fine-tuned on the COCO subset
- **Caption Generation:** BLIP generated descriptions for both real and synthetic images by interpreting detected object regions

Part B: Reinforcement Learning Agent (Theoretical Design)

Note: Part B is a theoretical contribution and was not implemented in the accompanying code.

- **Environment:** A Flappy Bird-style game recreated using OpenAI Gym interfaces
- **State Variables:** Vertical position and velocity of the bird, horizontal distance to the next pipe, and vertical position of the pipe gap
- **Action Space:**
 1. Flap (ascend)
 2. Do nothing (descend)
- **Network Architecture:**
 - Six input neurons
 - Two hidden layers with 128 and 64 neurons (ReLU activations)
 - Two output neurons representing Q-values for the available actions
- **Training Methodology:**
 - Epsilon-greedy policy for balancing exploration and exploitation
 - Experience replay with a buffer of 100,000 transitions
 - Target Q-network updated every 10,000 steps
- **Reward Structure:**
 - +1 for successfully passing through a set of pipes

- **-1** for collision with a pipe, the ground, or the ceiling
- **-0.01** as a small time penalty applied at each time step to encourage efficient flight

This reward system encourages the agent to prioritize survival, penalizes failure, and provides a dense feedback signal to accelerate convergence.

5. Analysis

Part A: Results and Observations

- The cGAN generated plausible class-conditional images, though visual fidelity decreased with increased scene complexity
- Faster R-CNN achieved 0.7616 mAP at an IoU threshold of 0.5, demonstrating strong object detection performance
- BLIP generated relevant captions for single-object and simple images but occasionally struggled with multi-object or occluded scenes

Part B: Theoretical Expectations

- The reward system (with +1, -1, and -0.01) provides dense, goal-aligned feedback
 - The time penalty discourages idle hovering, while the terminal rewards reinforce success and failure
 - Based on similar DQN architectures, convergence to a stable policy is expected within 1.5 to 2 million frames
 - Epsilon decay (e.g., exponential or linear) would support a gradual shift from exploration to exploitation
-

6. Summary and Conclusions

This report illustrates the dual strengths of deep learning: perception and action. The implemented vision pipeline demonstrates how generative and discriminative models can be combined to create a full-stack visual understanding system. The proposed DQN agent outlines a clear strategy for learning efficient gameplay policies using reinforcement learning.

While only Part A was implemented, the conceptual design of Part B aligns with modern reinforcement learning practices and is suitable for future implementation.

Future enhancements could include StyleGAN or diffusion models for higher-fidelity image generation, transformer-based object detection, or policy-gradient methods like PPO for reinforcement learning.

7. References

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2672–2680.

- Li, J., Selvaraju, R. R., Gotmare, A., Joty, S., Xiong, C., & Hoi, S. C. (2022). BLIP: Bootstrapping language-image pretraining for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., & others. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533.
- Mirza, M., & Osindero, S. (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
- Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6), 1137–1149.
-