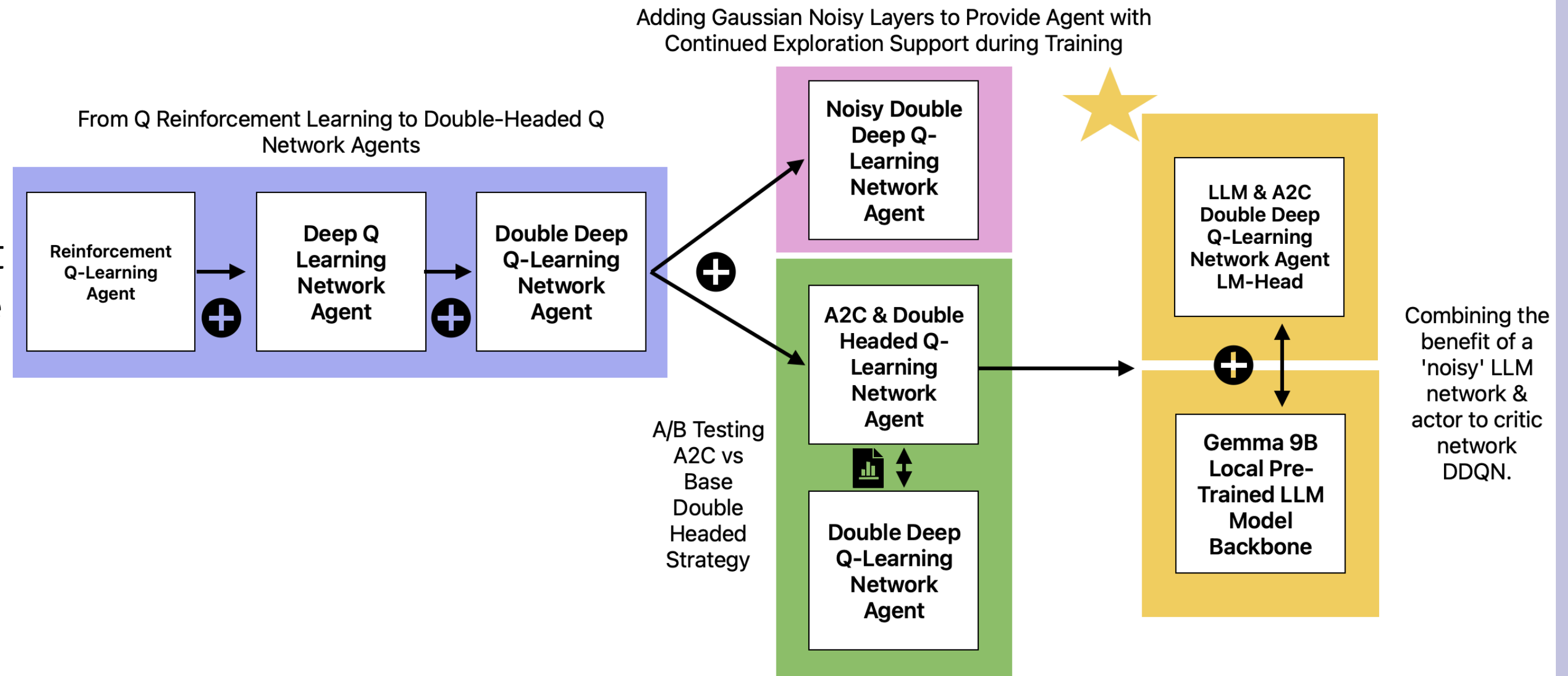# Path to Training a LLM-Based Deep Actor to Critic (A2C) & Deep Double Headed Q-Learning (DDQN) Neural Network Agent to Play Blackjack!
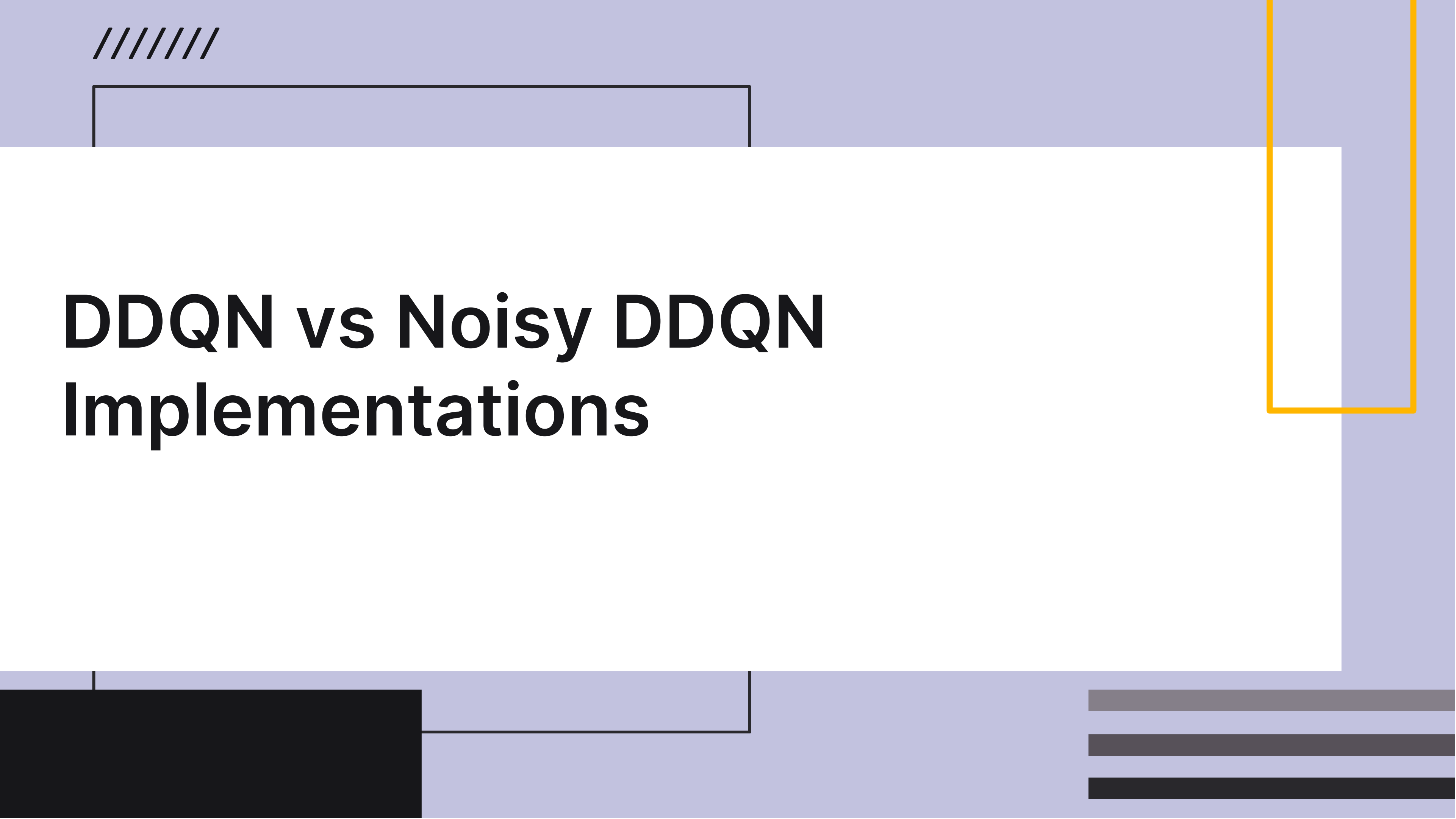


A project presentation by Alexis Ambriz, Sri Vinay Appari, Suresh Ravuri, and Harshavardhan Valmiki

# The Blackjack Agent Neural Network Architectures

As we embarked on the A.I Blackjack training journey, we had to remember that the key is to enjoy the game! The proposed A.I should learn the right mix of strategy and a "fun-loving" attitude with experimental playstyles.

Adding Gaussian Noisy Layers to Provide Agent with Continued Exploration Support during Training

From Q Reinforcement Learning to Double-Headed Q Network Agents

Reinforcement Q-Learning Agent **+** Deep Q Learning Network Agent **+** Double Deep Q-Learning Network Agent

**+**

Noisy Double Deep Q-Learning Network Agent

A/B Testing A2C vs Base Double Headed Strategy

A2C & Double Headed Q-Learning Network Agent

Double Deep Q-Learning Network Agent

LLM & A2C Double Deep Q-Learning Network Agent LM-Head

**+**

Gemma 9B Local Pre-Trained LLM Model Backbone

Combining the benefit of a 'noisy' LLM network & actor to critic network DDQN.

# DDQN vs Noisy DDQN Implementations

# DDQN Agent

```python
class ImprovedDQNAgent:
    def __init__(self, input_dim=3, learning_rate=5e-4, gamma=0.99, epsilon=1.0):
        self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
        self.fitness = [] # used to store a series of 'avg reward' during evaluation

        # Larger network
        self.policy_net = nn.Sequential(
            nn.Linear(input_dim, 128),
            nn.ReLU(),
            nn.Linear(128, 128),
            nn.ReLU(),
            nn.Linear(128, 64),
            nn.ReLU(),
            nn.Linear(64, 2)
        ).to(self.device)

        self.target_net = nn.Sequential(
            nn.Linear(input_dim, 128),
            nn.ReLU(),
            nn.Linear(128, 128),
            nn.ReLU(),
            nn.Linear(128, 64),
            nn.ReLU(),
            nn.Linear(64, 2)
        ).to(self.device)

        self.target_net.load_state_dict(self.policy_net.state_dict())

        self.optimizer = optim.Adam(self.policy_net.parameters(), lr=learning_rate)
        self.memory = ReplayMemory(capacity=20000)
```

## Predicts w/ Policy Net

```python
42    def select_action(self, state):
43        """Select action using epsilon-greedy policy"""
44        if random.random() < self.epsilon:
45            return random.randint(0, 1)
46
47        with torch.no_grad():
48            state = torch.FloatTensor(state).unsqueeze(0).to(self.device)
49            q_values = self.policy_net(state)
50            return q_values.argmax().item()
```

## Learns Q w/ Target Net

```python
# Double DQN implementation
with torch.no_grad():
    next_actions = self.policy_net(next_states).argmax(1)
    next_q_values = self.target_net(next_states).gather(1, next_actions.unsqueeze(1
    target_q_values = rewards + (1 - dones.float()) * self.gamma * next_q_values

current_q_values = self.policy_net(states).gather(1, actions.unsqueeze(1)).squeeze(

# Huber loss for better stability
loss = nn.SmoothL1Loss()(current_q_values, target_q_values)

self.optimizer.zero_grad()
loss.backward()
torch.nn.utils.clip_grad_norm_(self.policy_net.parameters(), 1.0)
self.optimizer.step()

return loss.item()
```

# Noisy DDQN Agent

## Noisy DDQN Architecture

```python
1  import torch.optim as optim
2
3  class NoisyDoubleDQNAgent:
4      def __init__(self, input_dim=3, learning_rate=5e-4, gamma=0.99, epsilon=1.0):
5          self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
6          self.fitness = [] # used to store a series of 'avg reward' during evaluation
7
8          self.policy_noisy_linear1 = NoisyLinear(128, 64)
9          self.policy_noisy_linear2 = NoisyLinear(64, 2)
10         self.target_noisy_linear1 = NoisyLinear(128, 64)
11         self.target_noisy_linear2 = NoisyLinear(64, 2)
12
13         # Larger network
14         self.policy_net = nn.Sequential(
15             nn.Linear(input_dim, 128),
16             nn.ReLU(),
17             nn.Linear(128, 128),
18             nn.ReLU(),
19             self.policy_noisy_linear1, # the noisy layers w/o ReLU in-between
20             self.policy_noisy_linear2,
21         ).to(self.device)
22
23         self.target_net = nn.Sequential(
24             nn.Linear(input_dim, 128),
25             nn.ReLU(),
26             nn.Linear(128, 128),
27             nn.ReLU(),
28             self.target_noisy_linear1,
29             self.target_noisy_linear2
30         ).to(self.device)
31
```

## Learns Gaussian Noise

```python
1  class NoisyLinear(nn.Module):
2      """Noisy linear module for NoisyNet.
3
4      Attributes:
5          in_features (int): input size of linear module
6          out_features (int): output size of linear module
7          std_init (float): initial std value
8          weight_mu (nn.Parameter): mean value weight parameter
9          weight_sigma (nn.Parameter): std value weight parameter
10         bias_mu (nn.Parameter): mean value bias parameter
11         bias_sigma (nn.Parameter): std value bias parameter
```

## Noise Added to Both Nets

```python
59     def forward(self, x: torch.Tensor) -> torch.Tensor:
60         """Forward method implementation.
61
62         We don't use separate statements on train / eval mode.
63         It doesn't show remarkable difference of performance.
64         """
65         return F.linear(
66             x,
67             self.weight_mu + self.weight_sigma * self.weight_epsilon,
68             self.bias_mu + self.bias_sigma * self.bias_epsilon,
```

# DDQN vs Actor-Critic (A2C) Implementation

## Similarity

The noisy DDQN / Q-Learning and off-policy agent was similar in architecture to an A2C agent as it also implemented separate neural network 'heads' to separate the target net (for training) from the policy net (for prediction), while also learning to add noise to the predictions of both. Furthermore, both the state and target policy are used to estimate optimal action values. Similarly, the replay buffer is used to train the policy net an episode behind the target network. The target network then copies the policy network during the next episode for prediction of actions.

## Difference

The A2C agent is a combo of SARSA and Q-Learning and is on-policy; architecture's key difference to the DDQN is that it uses a third network providing more separation. There is a feature extraction network, an 'actor' or action network, and a "critic" or a state / Q value network. Another difference is that all three networks are used during both training & prediction.

# A2C Agent

## Actor-Critic Network Architecture

```python
1  class ActorCritic(nn.Module):
2      def __init__(self, input_dim=3):
3          super(ActorCritic, self).__init__()
4
5          # Shared features extractor
6          self.features = nn.Sequential(
7              nn.Linear(input_dim, 128),
8              nn.ReLU(),
9              nn.Linear(128, 128),
10             nn.ReLU()
11         )
12
13         # Actor head (policy network)
14         self.actor = nn.Sequential(
15             nn.Linear(128, 64),
16             nn.ReLU(),
17             nn.Linear(64, 2),  # 2 actions: hit or stand
18             nn.Softmax(dim=-1)  # Output probabilities for each action
19         )
20
21         # Critic head (value network)
22         self.critic = nn.Sequential(
23             nn.Linear(128, 64),
24             nn.ReLU(),
25             nn.Linear(64, 1)  # Single value output
26         )
27
28     def forward(self, x):
29         features = self.features(x)
30         action_probs = self.actor(features)
31         state_value = self.critic(features)
32         return action_probs, state_value
```

```python
33  class A2CAgent:
34      def __init__(self, model, optimizer, gamma=0.99):
35          self.model = model
36          self.optimizer = optimizer
37          self.gamma = gamma
38          self.rewards = []
39          self.log_probs = []
40          self.state_values = []
41          self.entropies = []
42          self.device = torch.device("cuda" if torch.cuda.is_availab
43
44      def select_action(self, state):
45          """Select action using the policy network"""
46          state = torch.FloatTensor(state).to(self.device)
47          state.requires_grad = True
48
49          # Get action probabilities and state value
50          action_probs, state_value = self.model(state.unsqueeze(0))
51          action_probs = action_probs.squeeze()
52
53          # Create categorical distribution
54          dist = torch.distributions.Categorical(action_probs)
55          action = dist.sample()
56
57          # Calculate log probability and entropy
58          log_prob = dist.log_prob(action)
59          entropy = dist.entropy()
60
61          # Store for training
62          self.log_probs.append(log_prob)
63          self.state_values.append(state_value)
64          self.entropies.append(entropy)
65
```
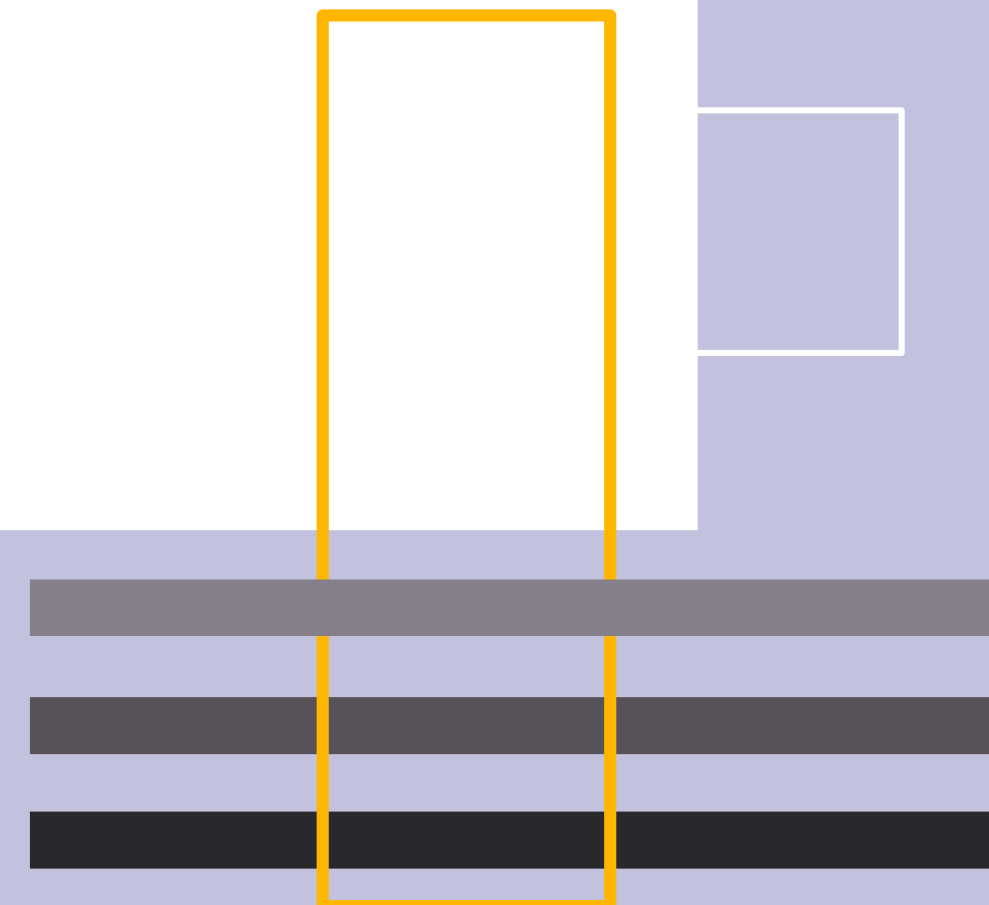
## A2C Agent vs the DDQN at Blackjack Strategy

When compared to the previous DDQN agent, the A2C agent demonstrated better performance. Across 200 test rounds:

- A2C Win Rate: 35.0%
- DQN Win Rate: 18.0%

During training they also differed:

- Average A2C Reward: -0.250
- Average DDQN Reward: -0.591
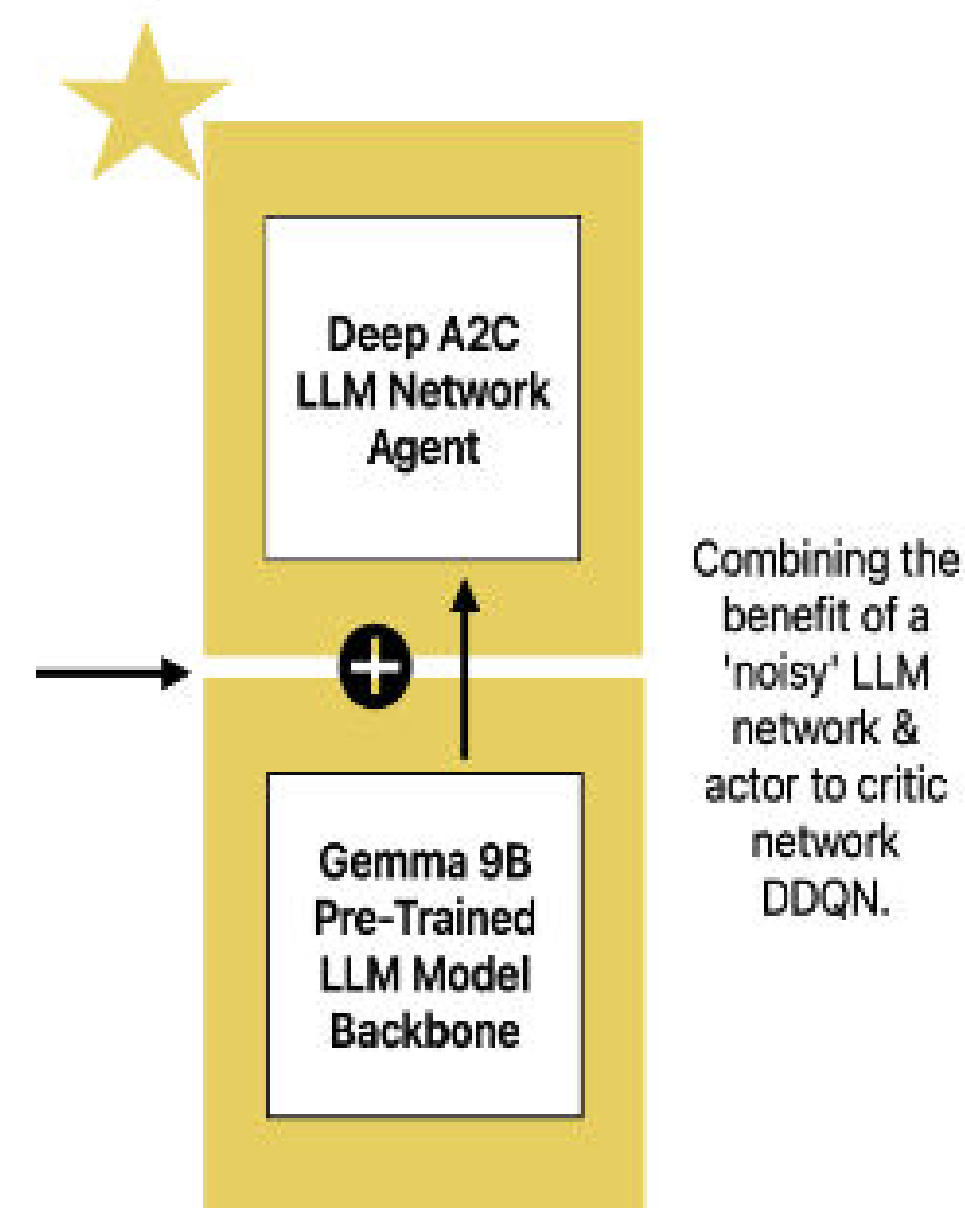- Decision Disagreement Rate: 48.0%

# Future Work

## A2C Agent LM-Head with Gemma 9B LLM Backbone

After evaluating the A2C reinforcement learning neural network agent and finding that it outperformed our previous DQN, future work could extend an A2C agent by leveraging the large amount of knowledge stored within popular foundation language models such as Gemma 9B or other larger models. Combining the A2C agent as the output 'head' with the language backbone could help diversify the agents learning strategies to potentially boost performance.

Deep A2C LLM Network Agent

Gemma 9B Pre-Trained LLM Model Backbone

Combining the benefit of a 'noisy' LLM network & actor to critic network DDQN.

# Thanks!