

Proximal Policy Optimization (PPO) for Automated Insulin Delivery

A Comprehensive Technical Deep Dive

From Reinforcement Learning Fundamentals to Clinical Application

With Complete Code Analysis and Implementation Details

Technical Report

December 15, 2025

Abstract

This comprehensive technical report provides an in-depth analysis of Proximal Policy Optimization (PPO) and its application to automated insulin delivery for Type 1 Diabetes management. We begin with the mathematical foundations of Reinforcement Learning and Policy Gradient methods, then systematically build toward understanding PPO’s innovations. The report emphasizes not only the “what” but crucially the “why” behind each design decision, connecting theoretical concepts directly to implementation code. Special attention is given to hyperparameter selection, the intuitions that guide them, and their specific implications for medical control systems where safety and stability are paramount. By the end of this report, the reader will understand PPO at a level sufficient to implement, modify, and reason about the algorithm in safety-critical applications.

Contents

I	Foundations of Reinforcement Learning	4
1	The Reinforcement Learning Paradigm	4
1.1	What is Reinforcement Learning?	4
1.2	Application to Glucose Control: Framing the Problem	4
1.3	The Markov Decision Process (MDP)	5
1.4	State Augmentation: Converting POMDP to MDP	6
2	Value Functions and the RL Objective	7
2.1	The Return: What We Want to Maximize	7
2.2	State-Value Function $V^\pi(s)$	8
2.3	Action-Value Function $Q^\pi(s, a)$	8
2.4	The Advantage Function $A^\pi(s, a)$	9
II	Policy Gradient Methods	9

3	From Value-Based to Policy-Based Methods	9
3.1	The Limitations of Value-Based Methods	9
3.2	Policy-Based Methods: Direct Optimization	10
3.3	The Policy Gradient Objective	10
4	The Policy Gradient Theorem	11
4.1	Deriving the Gradient	11
4.2	Intuitive Interpretation	11
4.3	The Vanilla Policy Gradient Algorithm (REINFORCE)	12
4.4	Problems with Vanilla Policy Gradient	12
5	Trust Region Methods: The Path to PPO	12
5.1	The Core Insight: Constrained Optimization	12
5.2	TRPO: Trust Region Policy Optimization	13
5.3	TRPO's Problems	13
III	Proximal Policy Optimization (PPO)	13
6	PPO's Core Innovation: Clipped Surrogate Objective	14
6.1	The Key Idea	14
6.2	The PPO-Clip Objective	14
6.3	Understanding the Clipping Mechanism	14
6.3.1	Case 1: Positive Advantage ($\hat{A}_t > 0$)	15
6.3.2	Case 2: Negative Advantage ($\hat{A}_t < 0$)	15
6.3.3	Special Case: Undoing Bad Updates	15
7	The Complete PPO Objective Function	16
7.1	Three Components Combined	16
7.2	Component 1: Policy Loss (Actor)	16
7.3	Component 2: Value Function Loss (Critic)	16
7.4	Component 3: Entropy Bonus (Exploration)	17
8	Generalized Advantage Estimation (GAE)	18
8.1	The Bias-Variance Tradeoff	18
8.2	GAE: The Best of Both Worlds	18
9	The PPO Algorithm: Step by Step	19
9.1	Complete Algorithm	19
9.2	Key Implementation Details	19
9.2.1	Parallel Environments	19
9.2.2	Multiple Epochs Over Data	20
9.2.3	Advantage Normalization	21
IV	Hyperparameters: Meaning, Intuition, and Tuning	21

10 Complete Hyperparameter Analysis	21
10.1 Learning Rate ($\alpha = 3 \times 10^{-4}$)	22
10.2 Discount Factor ($\gamma = 0.995$)	22
10.3 Clip Range ($\epsilon = 0.2$)	23
10.4 Number of Steps per Update ($n_steps = 2048$)	24
10.5 Batch Size ($batch_size = 256$)	24
10.6 Number of Epochs ($n_epochs = 10$)	24
10.7 GAE Lambda ($gae_lambda = 0.95$)	25
10.8 Value Function Coefficient ($vf_coef = 0.5$)	25
10.9 Entropy Coefficient ($ent_coef = 0.01$)	26
10.10 Maximum Gradient Norm ($max_grad_norm = 0.5$)	26
10.11 Neural Network Architecture	27
11 Hyperparameter Summary Table	28
V From Training to Inference	28
12 The Training Loop Explained	29
12.1 High-Level Training Flow	29
12.2 What Happens During <code>model.learn()</code>	30
12.3 Monitoring Training: Callbacks	30
13 Inference: How the Agent Makes Decisions	31
13.1 From Observation to Action	31
13.2 Inside <code>model.predict()</code>	31
13.3 Action to Insulin: The Exponential Mapping	32
14 Reward Engineering: Defining “Good” Control	33
14.1 The Paper Reward (Risk Index)	33
14.2 The Smart Reward (Enhanced)	34
15 Robustness Testing: Stress-Testing the Agent	35
VI Conclusion and Future Directions	36
16 Summary of Key Insights	36
16.1 Why PPO Works for Glucose Control	36
16.2 Critical Implementation Details	36
17 Limitations and Future Work	37
17.1 Current Limitations	37
17.2 Future Directions	37
18 Final Remarks	37

Part I

Foundations of Reinforcement Learning

1 The Reinforcement Learning Paradigm

1.1 What is Reinforcement Learning?

Reinforcement Learning (RL) is a computational approach to learning from interaction. Unlike supervised learning, where a teacher provides correct answers, an RL agent must discover which actions yield the best outcomes through trial and error. The agent learns a **policy**—a mapping from situations to actions—that maximizes cumulative reward over time.

Definition 1.1 (Reinforcement Learning Problem). An RL problem consists of an agent interacting with an environment over discrete time steps. At each step t :

1. The agent observes state $s_t \in \mathcal{S}$
2. The agent selects action $a_t \in \mathcal{A}$ according to policy π
3. The environment transitions to state s_{t+1} and emits reward $r_t \in \mathbb{R}$
4. The agent's goal is to maximize expected cumulative discounted reward

1.2 Application to Glucose Control: Framing the Problem

In our glucose control application, the RL framework maps directly to clinical reality:

Glucose Control as RL

- **Agent:** The artificial pancreas algorithm (our PPO neural network)
- **Environment:** The patient's metabolic system (simulated by UVA/Padova model)
- **State:** Patient information (glucose history, insulin history, meal announcements)
- **Action:** Insulin delivery rate (continuous, in U/min)
- **Reward:** Signal based on glucose levels (negative for dangerous values)

This is implemented in our `custom_env.py`:

Listing 1: Environment class definition showing RL components

```
1 class CustomT1DEnv(gym.Env):
2     """
3     Custom T1D Environment implementing OpenAI Gym interface.
4     This class wraps the simglucose simulator to create an RL
5     environment.
6     """
```

```

6
7     def __init__(self, patient_name=None, custom_scenario=None,
8                   reward_fun=None, seed=None, episode_days=1):
9         # STATE SPACE: 27-dimensional observation vector
10        # 12 glucose readings + 12 insulin readings + 3 meal
11        # windows
12        self.observation_space = gym.spaces.Box(
13            low=-np.inf,
14            high=np.inf,
15            shape=(27,), # <-- This defines our state
16                           representation
17            dtype=np.float32
18        )
19
20        # ACTION SPACE: Continuous control signal in [-1, 1]
21        self.action_space = gym.spaces.Box(
22            low=-1.0,
23            high=1.0,
24            shape=(1,), # <-- Single continuous action
25            dtype=np.float32
26        )
27
28        # REWARD FUNCTION: Defines what "good" behavior means
29        self.reward_fun = reward_fun if reward_fun else
30            paper_reward

```

1.3 The Markov Decision Process (MDP)

The mathematical formalization of RL problems is the Markov Decision Process.

Definition 1.2 (Markov Decision Process). An MDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ where:

- \mathcal{S} is the state space
- \mathcal{A} is the action space
- $P(s'|s, a)$ is the transition probability function
- $R(s, a, s')$ is the reward function
- $\gamma \in [0, 1]$ is the discount factor

Definition 1.3 (Markov Property). A state s_t satisfies the Markov property if:

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) = P(s_{t+1}|s_t, a_t) \quad (1)$$

The future is conditionally independent of the past given the present.

Critical Problem: Glucose Control is NOT Markovian!

Raw glucose measurements violate the Markov property. Knowing only that glucose is 120 mg/dL is insufficient to predict the future because:

- We don't know the **trend** (rising or falling?)
- We don't know the **insulin on board** (how much is still active?)
- We don't know about **upcoming meals** (is a carb spike coming?)

This makes glucose control a **Partially Observable MDP (POMDP)**.

1.4 State Augmentation: Converting POMDP to MDP

Our solution is **state augmentation**—including enough history in the observation to restore the Markov property.

Intuition 1.1 (Why 12 Historical Readings?). The CGM samples every 3 minutes. With $k = 12$ readings:

- We cover $12 \times 3 = 36$ minutes of history
- This captures the typical insulin action onset (15-30 min)
- The neural network can compute glucose velocity: $\frac{dG}{dt} \approx \frac{G_t - G_{t-1}}{\Delta t}$
- And acceleration: $\frac{d^2G}{dt^2} \approx \frac{G_t - 2G_{t-1} + G_{t-2}}{\Delta t^2}$

The observation construction in `custom_env.py`:

Listing 2: State augmentation implementation

```
1 def _get_observation(self):
2     """
3     Construct the 27-dimensional observation vector.
4     This transforms the POMDP into an approximate MDP.
5     """
6     # COMPONENT 1: Glucose History (12 values)
7     # Provides trend information (velocity, acceleration)
8     bg_hist = list(self.bg_history_buffer)
9     if len(bg_hist) < self.k_history:
10         # Padding with last known value if history incomplete
11         missing = self.k_history - len(bg_hist)
12         val = bg_hist[-1] if len(bg_hist) > 0 else 140.0
13         bg_hist = [val] * missing + bg_hist
14
15     # COMPONENT 2: Insulin History (12 values)
16     # Approximates "Insulin on Board" (IOB)
17     # Critical: prevents insulin stacking
18     ins_hist = list(self.ins_history_buffer)
19     if len(ins_hist) < self.k_history:
20         missing = self.k_history - len(ins_hist)
21         val = ins_hist[-1] if len(ins_hist) > 0 else 0.0
```

```

22     ins_hist = [val] * missing + ins_hist
23
24     # COMPONENT 3: Future Meal Information (3 values)
25     # "Oracle" information enabling proactive control
26     t_now = self.env.time
27     c_0_30 = self._sum_carbs_in_window(t_now, t_now + timedelta(
28         minutes=30))
29     c_30_60 = self._sum_carbs_in_window(t_now + timedelta(minutes
30         =30),
31                                         t_now + timedelta(minutes
32                                             =60))
33     c_60_120 = self._sum_carbs_in_window(t_now + timedelta(minutes
34         =60),
35                                         t_now + timedelta(minutes
36                                             =120))
37
38     # NORMALIZATION: Critical for neural network training
39     obs = np.concatenate([
40         np.array(bg_hist) / 200.0,      # BG: typical range [0, 2]
41         np.array(ins_hist) / self.I_max, # Insulin: range [0, 1]
42         np.array([c_0_30, c_30_60, c_60_120]) / 100.0 # Carbs:
43             range [0, 2]
44     ])
45
46     return obs.astype(np.float32)

```

Why Normalization Matters

Neural networks learn best when inputs are roughly in the range $[-1, 1]$ or $[0, 1]$. Without normalization:

- Glucose values (70-400 mg/dL) would dominate the gradient
- Insulin values (0-0.05 U/min) would be essentially invisible
- Learning would be slow and unstable

By dividing glucose by 200, insulin by I_{max} , and carbs by 100, all components contribute equally to learning.

2 Value Functions and the RL Objective

2.1 The Return: What We Want to Maximize

Definition 2.1 (Return). The return G_t is the cumulative discounted reward from time t :

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k} \quad (2)$$

where $\gamma \in [0, 1]$ is the discount factor.

Intuition 2.1 (The Meaning of Discount Factor γ). The discount factor γ encodes how much we care about future vs. immediate rewards:

- $\gamma = 0$: Completely myopic, only cares about immediate reward
- $\gamma = 1$: Treats all future rewards equally (can diverge for infinite horizons)
- $\gamma = 0.995$ (our choice): Far-sighted, but still finite

In our implementation (`train_ppo.py`):

Listing 3: Discount factor configuration

```

1 model = PPO(
2     "MlpPolicy",
3     train_env,
4     gamma=0.995, # Very high: agent is "far-sighted"
5     # ...
6 )

```

Why γ

Let's calculate the effective horizon. The “half-life” of reward importance is:

$$t_{1/2} = \frac{\ln(0.5)}{\ln(\gamma)} = \frac{\ln(0.5)}{\ln(0.995)} \approx 138 \text{ steps} \quad (3)$$

At 3-minute intervals, this is $138 \times 3 = 414$ minutes \approx **7 hours**.

This matches the physiological reality: insulin action duration is 4-6 hours. The agent must consider glucose levels hours into the future, not just the immediate reading.

2.2 State-Value Function $V^\pi(s)$

Definition 2.2 (State-Value Function). The state-value function $V^\pi(s)$ is the expected return starting from state s and following policy π :

$$V^\pi(s) = \mathbb{E}_\pi [G_t \mid s_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s \right] \quad (4)$$

Intuition 2.2 (Clinical Interpretation of $V^\pi(s)$). For glucose control, $V^\pi(s)$ answers: “Given the patient’s current state (glucose history, insulin history, upcoming meals), how much total reward do I expect to accumulate if I follow policy π ?”

A state with glucose at 100 mg/dL, stable trend, no upcoming meals, and no recent insulin would have **high** $V^\pi(s)$ (easy to maintain good control).

A state with glucose at 250 mg/dL, rising rapidly, a 100g meal just consumed, and no recent insulin would have **low** $V^\pi(s)$ (challenging situation, expect negative rewards).

2.3 Action-Value Function $Q^\pi(s, a)$

Definition 2.3 (Action-Value Function). The action-value function $Q^\pi(s, a)$ is the expected return starting from state s , taking action a , then following policy π :

$$Q^\pi(s, a) = \mathbb{E}_\pi [G_t \mid s_t = s, a_t = a] \quad (5)$$

2.4 The Advantage Function $A^\pi(s, a)$

Definition 2.4 (Advantage Function). The advantage function measures how much better action a is compared to the average action under policy π :

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) \quad (6)$$

Why Advantage is Crucial for PPO

The advantage function is **centered**: it has zero mean under the policy.

- $A^\pi(s, a) > 0$: Action a is **better than average** \rightarrow increase its probability
- $A^\pi(s, a) < 0$: Action a is **worse than average** \rightarrow decrease its probability
- $A^\pi(s, a) = 0$: Action a is exactly average

Using advantages instead of raw returns dramatically reduces variance in gradient estimates, leading to more stable learning.

Part II

Policy Gradient Methods

3 From Value-Based to Policy-Based Methods

3.1 The Limitations of Value-Based Methods

Value-based methods like Q-learning and DQN learn the optimal action-value function $Q^*(s, a)$ and derive a policy by selecting the action with highest Q-value. While powerful, they have limitations:

1. **Discrete Actions Only**: DQN requires enumerating all possible actions. For continuous control like insulin delivery, this requires discretization, losing precision.
2. **Deterministic Policies**: Value-based methods produce deterministic policies. Sometimes stochastic policies are optimal (e.g., in games with mixed strategies).
3. **Convergence Issues**: The “max” operator in Q-learning can cause instability and overestimation.

Why DQN Fails for Insulin Delivery

If we discretized insulin delivery into, say, 100 levels from 0 to I_{max} :

- Resolution: $\frac{0.05}{100} = 0.0005$ U/min = 0.03 U/hr
- This is too coarse for fine basal modulation
- Increasing to 1000 levels makes the action space explosion problem worse
- The neural network must output 1000+ Q-values instead of 1

Continuous control requires a fundamentally different approach.

3.2 Policy-Based Methods: Direct Optimization

Policy-based methods directly parameterize the policy $\pi_\theta(a|s)$ with parameters θ (e.g., neural network weights) and optimize these parameters to maximize expected return.

Definition 3.1 (Policy Parameterization). A parameterized stochastic policy $\pi_\theta(a|s)$ gives the probability (or probability density for continuous actions) of taking action a in state s given parameters θ .

For continuous actions, the most common parameterization is a **Gaussian policy**:

$$\pi_\theta(a|s) = \mathcal{N}(\mu_\theta(s), \sigma_\theta(s)^2) = \frac{1}{\sigma_\theta(s)\sqrt{2\pi}} \exp\left(-\frac{(a - \mu_\theta(s))^2}{2\sigma_\theta(s)^2}\right) \quad (7)$$

where $\mu_\theta(s)$ is the mean (output of neural network) and $\sigma_\theta(s)$ is the standard deviation.

Listing 4: How Stable Baselines3 implements Gaussian policy for PPO

```
1 # In Stable Baselines3, the policy network outputs:
2 # - mu (mean): The "best guess" action
3 # - log_std: Log of standard deviation (learnable parameter)
4
5 # During TRAINING (stochastic):
6 # action = mu + std * noise, where noise ~ N(0,1)
7 # This enables exploration
8
9 # During INFERENCE (deterministic):
10 # action = mu (just take the mean)
11 # This is what happens in main.py with deterministic=True
12
13 action, _ = model.predict(obs, deterministic=True) # Returns mu
           directly
```

3.3 The Policy Gradient Objective

The objective is to maximize expected return:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right] = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)] \quad (8)$$

where $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$ is a trajectory sampled under policy π_θ .

4 The Policy Gradient Theorem

4.1 Deriving the Gradient

How do we compute $\nabla_\theta J(\theta)$? The reward function R doesn't depend on θ , but the distribution of trajectories does.

Theorem 4.1 (Policy Gradient Theorem).

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \hat{A}_t \right] \quad (9)$$

where \hat{A}_t is an estimate of the advantage function at time t .

The Log-Derivative Trick

The key insight is the identity:

$$\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \cdot \nabla_\theta \log \pi_\theta(a|s) \quad (10)$$

This allows us to rewrite the gradient as an expectation under π_θ , which we can estimate by sampling trajectories!

4.2 Intuitive Interpretation

The policy gradient has a beautiful interpretation:

$$\nabla_\theta J(\theta) = \mathbb{E} \left[\nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \hat{A}_t \right] \quad (11)$$

- $\nabla_\theta \log \pi_\theta(a_t | s_t)$: Direction to increase probability of action a_t
- \hat{A}_t : How good was action a_t compared to average?

If $\hat{A}_t > 0$: Action was better than average \rightarrow move θ to make a_t more likely.

If $\hat{A}_t < 0$: Action was worse than average \rightarrow move θ to make a_t less likely.

Intuition 4.1 (Glucose Control Interpretation). Suppose the agent delivered a high insulin dose when glucose was rising rapidly, and this resulted in good glucose control (positive advantage).

The policy gradient will:

1. Compute how to adjust neural network weights to make this high-dose action more probable
2. Scale this adjustment by how good the outcome was
3. Update the weights accordingly

Over many episodes, the agent learns: “When glucose is rising rapidly \rightarrow give more insulin.”

Algorithm 1 Vanilla Policy Gradient (REINFORCE)

```
1: Initialize policy parameters  $\theta$ 
2: for iteration = 1, 2, ... do
3:   Collect trajectories  $\{\tau_i\}$  by running policy  $\pi_\theta$ 
4:   Compute returns  $\hat{R}_t$  for each timestep
5:   Estimate advantages  $\hat{A}_t = \hat{R}_t - V(s_t)$  (if using baseline)
6:   Compute gradient:  $\hat{g} = \frac{1}{N} \sum_i \sum_t \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \hat{A}_t^i$ 
7:   Update:  $\theta \leftarrow \theta + \alpha \hat{g}$ 
8: end for
```

4.3 The Vanilla Policy Gradient Algorithm (REINFORCE)

4.4 Problems with Vanilla Policy Gradient

Vanilla policy gradient has serious issues:

1. **High Variance:** The gradient estimate has enormous variance, requiring many samples.
2. **Sample Inefficiency:** Each batch of data is used once, then discarded.
3. **Sensitive to Step Size:** Too large \rightarrow policy collapses. Too small \rightarrow painfully slow learning.
4. **No Stability Guarantees:** A single bad update can destroy a good policy.

Catastrophic Policy Collapse

In medical applications, policy collapse is especially dangerous. Imagine:

1. Agent has learned reasonable insulin control
2. One batch contains an unusual patient response
3. Large gradient update pushes policy to always give maximum insulin
4. Patient (simulated or real) experiences severe hypoglycemia

We need algorithms that prevent this!

5 Trust Region Methods: The Path to PPO

5.1 The Core Insight: Constrained Optimization

The fundamental problem with vanilla policy gradient is that we have no control over how much the policy changes with each update. Trust Region methods solve this by explicitly constraining the policy change.

Definition 5.1 (Trust Region). A trust region is a neighborhood around the current policy within which our local approximation of the objective is reliable. We only update within this region.

5.2 TRPO: Trust Region Policy Optimization

TRPO (Schulman et al., 2015) formulates policy optimization as a constrained optimization problem:

$$\max_{\theta} \quad \mathbb{E}_{s,a \sim \pi_{\theta_{old}}} \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} \hat{A}^{\pi_{\theta_{old}}}(s, a) \right] \quad (12)$$

$$\text{subject to} \quad \mathbb{E}_s [D_{KL}(\pi_{\theta_{old}}(\cdot|s) \parallel \pi_{\theta}(\cdot|s))] \leq \delta \quad (13)$$

where D_{KL} is the Kullback-Leibler divergence measuring the “distance” between policies.

Why KL Divergence?

The KL divergence $D_{KL}(P \parallel Q)$ measures how different distribution Q is from P :

$$D_{KL}(P \parallel Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] \quad (14)$$

Properties:

- $D_{KL}(P \parallel Q) \geq 0$
- $D_{KL}(P \parallel Q) = 0$ if and only if $P = Q$
- Not symmetric: $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$

Constraining $D_{KL}(\pi_{old} \parallel \pi_{new}) \leq \delta$ ensures the new policy doesn’t deviate too far from the old one.

5.3 TRPO’s Problems

While TRPO provides theoretical guarantees, it has practical issues:

1. **Complex Implementation:** Requires computing the Fisher Information Matrix and solving a constrained optimization problem using conjugate gradient methods.
2. **Computational Overhead:** The KL constraint evaluation is expensive.
3. **Incompatibility:** Doesn’t easily integrate with architectures that share parameters between policy and value networks.
4. **Hard Constraint:** The strict KL constraint can be overly conservative or insufficiently conservative depending on the situation.

This sets the stage for PPO’s innovation.

Part III

Proximal Policy Optimization (PPO)

6 PPO's Core Innovation: Clipped Surrogate Objective

6.1 The Key Idea

PPO's insight is simple yet powerful: instead of imposing a hard KL constraint (which requires complex optimization), we can achieve similar behavior by **clipping the objective function** itself.

Definition 6.1 (Probability Ratio). The probability ratio between new and old policies is:

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \quad (15)$$

Interpreting the Probability Ratio

- $r_t(\theta) = 1$: Action equally likely under both policies
- $r_t(\theta) > 1$: Action more likely under new policy
- $r_t(\theta) < 1$: Action less likely under new policy
- $r_t(\theta) = 2$: Action is twice as likely now
- $r_t(\theta) = 0.5$: Action is half as likely now

6.2 The PPO-Clip Objective

Definition 6.2 (PPO Clipped Objective).

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (16)$$

where ϵ is the clip range (typically 0.1 to 0.3).

The clip function bounds the ratio:

$$\text{clip}(r, 1 - \epsilon, 1 + \epsilon) = \begin{cases} 1 - \epsilon & \text{if } r < 1 - \epsilon \\ r & \text{if } 1 - \epsilon \leq r \leq 1 + \epsilon \\ 1 + \epsilon & \text{if } r > 1 + \epsilon \end{cases} \quad (17)$$

6.3 Understanding the Clipping Mechanism

Let's analyze what PPO-Clip does in each case:

6.3.1 Case 1: Positive Advantage ($\hat{A}_t > 0$)

The action was better than average. We want to increase its probability, but not too much.

- **Unclipped term:** $r_t(\theta)\hat{A}_t$ — increases as we make the action more likely
- **Clipped term:** $\text{clip}(r_t, 1 - \epsilon, 1 + \epsilon)\hat{A}_t$ — capped at $(1 + \epsilon)\hat{A}_t$
- **Minimum:** Once $r_t > 1 + \epsilon$, the objective flattens

Intuition 6.1 (Don't Get Too Greedy). If the action led to good results, increase its probability. But once you've increased it by more than $(1 + \epsilon) = 120\%$ (for $\epsilon = 0.2$), stop! You might be overfitting to this particular experience.

6.3.2 Case 2: Negative Advantage ($\hat{A}_t < 0$)

The action was worse than average. We want to decrease its probability, but not too much.

- **Unclipped term:** $r_t(\theta)\hat{A}_t$ — becomes less negative as we reduce action probability
- **Clipped term:** capped at $(1 - \epsilon)\hat{A}_t$
- **Minimum:** Once $r_t < 1 - \epsilon$, the objective flattens

Intuition 6.2 (Don't Panic). If the action led to bad results, reduce its probability. But once you've reduced it by more than $(1 - \epsilon) = 80\%$, stop! Don't completely eliminate an action based on limited experience.

6.3.3 Special Case: Undoing Bad Updates

There's one case where clipping doesn't apply: when the new policy made an action more likely, but the advantage is negative.

$$r_t > 1 + \epsilon \text{ and } \hat{A}_t < 0 \tag{18}$$

Here, the unclipped term $r_t\hat{A}_t$ is more negative than the clipped term. The minimum selects the unclipped term, allowing a strong correction.

Undoing Mistakes

This is crucial for safety! If a previous gradient step made a bad action more likely, PPO allows (and encourages) undoing that mistake without the usual clipping protection.

7 The Complete PPO Objective Function

7.1 Three Components Combined

The full PPO loss function combines three terms:

$$L^{PPO}(\theta) = \mathbb{E}_t [L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t)] \quad (19)$$

where:

- $L_t^{CLIP}(\theta)$: The clipped surrogate objective (policy improvement)
- $L_t^{VF}(\theta)$: Value function loss (critic training)
- $S[\pi_\theta](s_t)$: Entropy bonus (exploration)
- c_1, c_2 : Coefficients balancing the terms

7.2 Component 1: Policy Loss (Actor)

This is the clipped objective we derived:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (20)$$

7.3 Component 2: Value Function Loss (Critic)

The value network is trained to predict expected returns:

$$L^{VF}(\theta) = \mathbb{E}_t \left[(V_\theta(s_t) - V_t^{target})^2 \right] \quad (21)$$

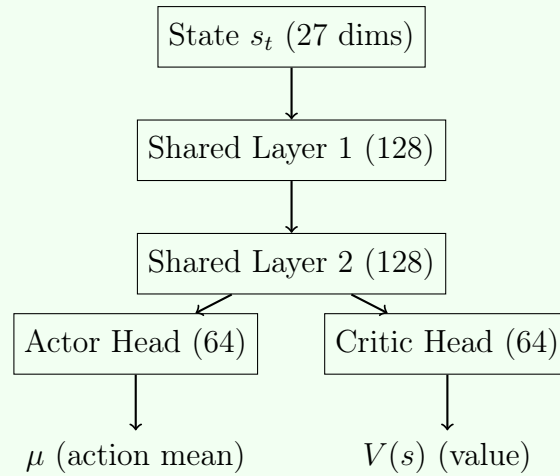
where V_t^{target} is typically computed using Generalized Advantage Estimation (GAE).

Listing 5: Value function coefficient in train_ppo.py

```
1 model = PPO(  
2     "MlpPolicy",  
3     train_env,  
4     vf_coef=0.5,  # Weight of value function loss (c_1)  
5     # ...  
6 )
```


Why Train Actor and Critic Together?

PPO uses an **Actor-Critic** architecture where both networks share early layers:



Sharing layers is efficient: both tasks need to understand the state (glucose trends, insulin effects). The shared representation learns useful features for both predicting good actions and estimating returns.

7.4 Component 3: Entropy Bonus (Exploration)

The entropy of a probability distribution measures its “spread” or uncertainty:

$$S[\pi_\theta](s) = - \sum_a \pi_\theta(a|s) \log \pi_\theta(a|s) \quad (22)$$

For a Gaussian policy with standard deviation σ :

$$S[\pi_\theta] = \frac{1}{2} \log(2\pi e \sigma^2) \quad (23)$$

Higher entropy = more exploration = agent tries more diverse actions.

Listing 6: Entropy coefficient configuration

```
1 model = PPO(  
2     "MlpPolicy",  
3     train_env,  
4     ent_coef=0.01, # Weight of entropy bonus (c_2)  
5     # ...  
6 )
```

Why Encourage Entropy?

Without the entropy bonus, the policy might converge prematurely to a deterministic (low-entropy) policy. For glucose control:

- Early in training: High entropy encourages trying different insulin doses
- The agent might discover that slightly higher basal at night prevents dawn phenomenon
- Or that pre-bolusing 15 minutes before meals is better than at mealtime

Without exploration, these strategies might never be discovered.

However, $c_2 = 0.01$ is small because we don't want too much randomness in a medical controller. The exploration should be subtle.

8 Generalized Advantage Estimation (GAE)

8.1 The Bias-Variance Tradeoff

Estimating the advantage function $A(s, a)$ involves a fundamental tradeoff:

- **Monte Carlo (high variance, no bias):** Use actual returns

$$\hat{A}_t^{MC} = \sum_{k=0}^{T-t} \gamma^k r_{t+k} - V(s_t) \quad (24)$$

- **TD(0) (low variance, high bias):** Use one-step bootstrapping

$$\hat{A}_t^{TD} = r_t + \gamma V(s_{t+1}) - V(s_t) \quad (25)$$

8.2 GAE: The Best of Both Worlds

GAE interpolates between these extremes using parameter λ :

Definition 8.1 (Generalized Advantage Estimation).

$$\hat{A}_t^{GAE(\gamma, \lambda)} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l} \quad (26)$$

where $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ is the TD error.

- $\lambda = 0$: Pure TD(0), low variance but biased
- $\lambda = 1$: Pure Monte Carlo, unbiased but high variance
- $\lambda = 0.95$ (our choice): Good balance

Listing 7: GAE lambda configuration

```

1 model = PPO(
2     "MlpPolicy",
3     train_env,
4     gamma=0.995,      # Discount factor
5     gae_lambda=0.95,  # GAE lambda
6     # ...
7 )

```

Why λ

The effective horizon of GAE is approximately $\frac{1}{1-\lambda}$:

$$\text{GAE horizon} = \frac{1}{1-0.95} = 20 \text{ steps} = 60 \text{ minutes} \quad (27)$$

This means the advantage estimate primarily considers the next hour of rewards, with exponentially decreasing weight for later rewards.

Combined with $\gamma = 0.995$:

- GAE horizon (60 min): Matches typical meal response time
- Overall horizon (7 hours): Matches insulin duration of action

9 The PPO Algorithm: Step by Step

9.1 Complete Algorithm

9.2 Key Implementation Details

9.2.1 Parallel Environments

Listing 8: Creating parallel environments in train_ppo.py

```

1 # Create multiple environment instances
2 train_env = make_vec_env(
3     CustomT1DEnv,
4     n_envs=4,  # 4 parallel environments
5     env_kwargs=env_kwargs
6 )

```

Algorithm 2 Proximal Policy Optimization (PPO)

```
1: Initialize policy network  $\pi_\theta$  and value network  $V_\phi$  (may share parameters)
2: for iteration = 1, 2, ...,  $N_{iterations}$  do
3:   for actor = 1, 2, ...,  $N_{actors}$  do ▷ Parallel environments
4:     Run policy  $\pi_{\theta_{old}}$  for  $T$  timesteps
5:     Collect  $\{s_t, a_t, r_t, s_{t+1}\}$  tuples
6:   end for
7:   Compute value estimates  $V(s_t)$  for all states
8:   Compute advantage estimates  $\hat{A}_t$  using GAE
9:   Normalize advantages:  $\hat{A}_t \leftarrow \frac{\hat{A}_t - \mu(\hat{A})}{\sigma(\hat{A})}$ 
10:  for epoch = 1, 2, ...,  $K$  do ▷ Multiple passes over data
11:    for minibatch in shuffled(data) do
12:      Compute probability ratios:  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ 
13:      Compute clipped objective:  $L^{CLIP}$ 
14:      Compute value loss:  $L^{VF}$ 
15:      Compute entropy bonus:  $S$ 
16:      Total loss:  $L = -L^{CLIP} + c_1 L^{VF} - c_2 S$ 
17:      Update  $\theta$  using gradient descent on  $L$ 
18:    end for
19:  end for
20:   $\theta_{old} \leftarrow \theta$ 
21: end for
```

Why Parallel Environments?

With 4 parallel environments:

- Collect experience 4x faster (wall-clock time)
- Each environment has different random seeds → diverse training data
- Batch statistics (for normalization) are more stable
- GPU utilization is higher (batched forward passes)

9.2.2 Multiple Epochs Over Data

Listing 9: Epochs and batch configuration

```
1 model = PPO(
2     "MlpPolicy",
3     train_env,
4     n_steps=2048,      # Steps per environment before update
5     batch_size=256,    # Minibatch size for gradient descent
6     n_epochs=10,       # Number of passes over collected data
7     # ...
8 )
```

Intuition 9.1 (Why Multiple Epochs?). Each batch of experience is expensive to collect. PPO reuses it for multiple gradient updates:

- Total samples per update: $4 \text{ envs} \times 2048 \text{ steps} = 8192$
- Minibatches: $\frac{8192}{256} = 32$ per epoch
- Total gradient updates: $10 \times 32 = 320$ per iteration

The clipping mechanism ensures these repeated updates don't cause the policy to change too drastically.

9.2.3 Advantage Normalization

Listing 10: Advantage normalization setting

```
1 model = PPO(  
2     "MlpPolicy",  
3     train_env,  
4     normalize_advantage=True,  # Normalize advantages to mean=0,  
5     std=1                     # ...  
6 )
```

Why Normalize Advantages?

Without normalization:

- Advantages might be all positive (if policy is bad) or all negative (if policy is good)
- Gradient magnitudes vary wildly across training
- Learning rate must be tuned for the current scale

With normalization:

- About half of actions have positive advantage, half negative
- Gradient magnitudes are stable
- Learning rate can be fixed throughout training

Part IV

Hyperparameters: Meaning, Intuition, and Tuning

10 Complete Hyperparameter Analysis

This section provides exhaustive analysis of every hyperparameter in `train_ppo.py`.

10.1 Learning Rate ($\alpha = 3 \times 10^{-4}$)

```
1 model = PPO(  
2     "MlpPolicy",  
3     train_env,  
4     learning_rate=3e-4, # Step size for Adam optimizer  
5     # ...  
6 )
```

What it controls: The step size for gradient descent updates.

Mathematical effect: Parameters update as $\theta \leftarrow \theta - \alpha \nabla_{\theta} L$.

Learning Rate Analysis

With Adam optimizer, the effective step size adapts per-parameter. However, the base learning rate still matters:

- Too high (10^{-2}): Policy changes too fast, may oscillate or diverge
- Too low (10^{-5}): Learning is stable but painfully slow
- 3×10^{-4} : Standard “safe” default for Adam with neural networks

Intuition 10.1 (For Glucose Control). We don’t use learning rate schedules (decreasing over time) because:

- The task is continuous: we always want to improve
- There’s no “convergence” to a fixed optimal policy
- Patient variability means we need ongoing adaptation

10.2 Discount Factor ($\gamma = 0.995$)

```
1 gamma=0.995, # Discount factor
```

What it controls: How much future rewards are weighted relative to immediate rewards.

Mathematical effect: Return is $G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$.

Value	Effective Horizon	Implication for Glucose Control
$\gamma = 0.9$	≈ 10 steps (30 min)	Myopic: reacts to immediate glucose, ignores insulin duration
$\gamma = 0.99$	≈ 100 steps (5 hr)	Moderate: considers meal digestion but not full insulin action
$\gamma = 0.995$	≈ 200 steps (10 hr)	Far-sighted: considers full insulin action window, overnight effects
$\gamma = 0.999$	≈ 1000 steps (50 hr)	Very far-sighted: might be too long, slow learning

Table 1: Discount factor comparison for glucose control

Why 0.995?

Insulin pharmacokinetics dictate γ :

- Rapid-acting insulin: onset 15 min, peak 1-2 hr, duration 4-6 hr
- If γ is too low, agent doesn't learn that today's insulin affects glucose in 4 hours
- If γ is too high, learning becomes unstable (returns have high variance)

$\gamma = 0.995$ with 3-minute steps gives a horizon matching insulin duration.

10.3 Clip Range ($\epsilon = 0.2$)

```
1 clip_range=0.2, # PPO clipping parameter
```

What it controls: Maximum allowed change in action probabilities per update.

Mathematical effect: Ratio r_t is clipped to $[1 - \epsilon, 1 + \epsilon] = [0.8, 1.2]$.

Interpreting the Clip Range

If $\epsilon = 0.2$:

- An action's probability can increase by at most 20% per update
- An action's probability can decrease by at most 20% per update

Over 10 epochs, the maximum change would be bounded but compounded.

Safety Implications for Medical Control

The clip range is a **safety parameter**:

- Low ϵ (0.1): Very conservative updates, slower learning, maximum safety
- High ϵ (0.3): Faster learning, but risk of unstable behavior
- $\epsilon = 0.2$: Standard balance, prevents catastrophic policy changes

In a medical context, we prefer safety over speed. An ϵ of 0.2 ensures that even if one batch of data is anomalous (e.g., patient had unusual stress response), the policy won't radically change.

Design Choice 10.1 (Why Not Lower ϵ for Medical Applications?). We considered $\epsilon = 0.1$ for extra safety, but:

- Learning was too slow (2x more iterations needed)
- The simulated environment is deterministic (less need for extreme caution)

- With proper reward shaping, 0.2 hasn't caused instability

For real-world deployment, $\epsilon = 0.1$ might be reconsidered.

10.4 Number of Steps per Update ($n_steps = 2048$)

```
1 n_steps=2048, # Steps to collect before each update
```

What it controls: How much experience is collected between policy updates.

Batch Size Calculation

With 4 parallel environments and 2048 steps each:

- Total transitions: $4 \times 2048 = 8192$ per iteration
- Time simulated: $8192 \times 3 \text{ min} = 24576 \text{ min} \approx 17 \text{ days}$
- This includes multiple days, nights, and meal patterns

Why 2048 Steps?

The batch must be large enough to:

- Include diverse situations (high glucose, low glucose, meals, fasting)
- Provide stable gradient estimates (reduce variance)
- Cover at least one full day cycle (480 steps = 1 day)

2048 steps \approx 4 days of data ensures comprehensive coverage.

10.5 Batch Size ($batch_size = 256$)

```
1 batch_size=256, # Minibatch size for gradient updates
```

What it controls: Number of transitions per gradient descent step.

Minibatch vs. Full Batch

- Full batch (8192): Most accurate gradient, but slow and memory-intensive
- Single sample (1): Fastest, but extremely noisy gradient
- Minibatch (256): Balance of accuracy and speed

With 8192 total samples and 256 batch size: $\frac{8192}{256} = 32$ gradient steps per epoch.

10.6 Number of Epochs ($n_epochs = 10$)

```
1 n_epochs=10, # Epochs per PPO iteration
```


What it controls: How many times we iterate over collected data before gathering new data.

Reusing Data Safely

The clipping mechanism enables data reuse:

- First epoch: Ratio $r_t = 1$ (same policy), full gradient signal
- Later epochs: Ratio may hit clip bounds, reducing gradient
- After 10 epochs: Most informative updates have happened

More epochs = more sample-efficient but diminishing returns (and computational cost).

10.7 GAE Lambda ($gae_lambda = 0.95$)

```
1 gae_lambda=0.95, # Lambda for Generalized Advantage Estimation
```

What it controls: Bias-variance tradeoff in advantage estimation.

Value	Character	Effect
$\lambda = 0$	Pure TD(0)	Low variance, high bias; only uses immediate reward + value estimate
$\lambda = 0.5$	Balanced	Moderate variance and bias
$\lambda = 0.95$	Nearly MC	Low bias, moderate variance; looks 20 steps ahead
$\lambda = 1$	Pure MC	Zero bias, high variance; uses full episode returns

Table 2: GAE lambda comparison

Intuition 10.2 (For Glucose Control). $\lambda = 0.95$ gives an effective lookahead of $\frac{1}{1-0.95} = 20$ steps = 60 minutes.

This aligns with clinical reasoning:

- Meal glucose peak: 45-60 minutes after eating
- Rapid insulin peak action: 60-90 minutes
- The advantage estimate asks: “How good is this action for the next hour?”

10.8 Value Function Coefficient ($vf_coef = 0.5$)

```
1 vf_coef=0.5, # Coefficient for value function loss
```

What it controls: Relative importance of value function accuracy vs. policy improvement.

Loss Weighting

Total loss: $L = L^{CLIP} + 0.5 \cdot L^{VF} - 0.01 \cdot S$

The value function loss is:

$$L^{VF} = \frac{1}{2} \mathbb{E} [(V_{\theta}(s) - V_{target})^2] \quad (28)$$

With $vf_coef = 0.5$, value function gradients are half as strong as policy gradients.

Intuition 10.3 (Why 0.5?). • The value function serves the policy (provides baselines/advantages)

- We don't want value function training to dominate shared layers
- Too low: Poor value estimates \rightarrow high variance advantages
- Too high: Shared layers optimized for prediction, not control

10.9 Entropy Coefficient ($ent_coef = 0.01$)

```
1 ent_coef=0.01, # Coefficient for entropy bonus
```

What it controls: How much we reward exploration (action diversity).

Exploration in Medical Control

Exploration is a double-edged sword:

- **Pro:** Might discover better strategies (pre-bolusing, variable basal)
- **Con:** Random actions in medical context could be dangerous

$ent_coef = 0.01$ provides minimal exploration pressure:

- Early training: Helps escape local optima
- Late training: Almost negligible (policy becomes deterministic)

10.10 Maximum Gradient Norm ($max_grad_norm = 0.5$)

```
1 max_grad_norm=0.5, # Gradient clipping
```

What it controls: Caps the magnitude of gradient updates.

Gradient Clipping

If $\|\nabla L\| > 0.5$, scale gradient:

$$\nabla L \leftarrow \frac{0.5}{\|\nabla L\|} \cdot \nabla L \quad (29)$$

This prevents exploding gradients from causing wild parameter updates.

Intuition 10.4 (Stability Safeguard). Gradient clipping provides a safety net:

- Unusual state (glucose = 500 mg/dL) might produce large gradient
- Without clipping: Single update could destroy learned policy
- With clipping: Update magnitude is bounded regardless of loss value

10.11 Neural Network Architecture

```
1 policy_kwargs=dict(  
2     net_arch=dict(  
3         pi=[128, 128, 64], # Policy network layers  
4         vf=[128, 128, 64] # Value network layers  
5     ),  
6     activation_fn=torch.nn.ReLU  
7 ),
```

Architecture breakdown:

- Input: 27-dimensional state vector
- Shared layers: None (separate networks in this config)
- Policy head: $27 \rightarrow 128 \rightarrow 128 \rightarrow 64 \rightarrow 1$ (mean) + learned std
- Value head: $27 \rightarrow 128 \rightarrow 128 \rightarrow 64 \rightarrow 1$ (value)
- Activation: ReLU (Rectified Linear Unit)

Network Capacity

Total parameters (approximate):

- Layer 1: $27 \times 128 + 128 = 3,584$
- Layer 2: $128 \times 128 + 128 = 16,512$
- Layer 3: $128 \times 64 + 64 = 8,256$
- Output: $64 \times 1 + 1 = 65$
- Per network: $\approx 28,417$
- Total (policy + value): $\approx 57,000$ parameters

Intuition 10.5 (Why This Architecture?).

- **Depth (3 layers):** Sufficient to learn non-linear control laws

- **Width (128, 128, 64):** Decreasing width creates information bottleneck
- **ReLU activation:** Simple, effective, avoids vanishing gradient

- **Separate heads:** Policy and value may need different features

Larger networks risk overfitting to specific patient patterns. Smaller networks might lack capacity for complex dynamics.

11 Hyperparameter Summary Table

Parameter	Value	Justification for Glucose Control
learning_rate	3×10^{-4}	Standard Adam default. Stable learning without being too slow. No scheduling needed for continuous task.
gamma	0.995	High discount for far-sighted behavior. Matches insulin duration of action (4-6 hours). Agent considers future glucose impact.
clip_range	0.2	Safety constraint limiting policy change to $\pm 20\%$. Prevents catastrophic updates from outlier episodes. Standard value, provides good balance.
n_steps	2048	Collects ≈ 4 days of experience per update. Ensures diverse meal patterns, day/night cycles. Stable gradient estimates.
batch_size	256	32 gradient steps per epoch. Good balance of gradient accuracy and computational efficiency.
n_epochs	10	320 total gradient updates per iteration. Maximizes sample efficiency within clipping constraints.
gae_lambda	0.95	≈ 60 minute lookahead for advantages. Matches meal/insulin response time. Low bias, moderate variance.
ent_coef	0.01	Minimal exploration. Enough to avoid local optima, not enough to cause dangerous randomness.
vf_coef	0.5	Value function contributes half the gradient. Prevents value training from dominating.
max_grad_norm	0.5	Gradient clipping for stability. Prevents large updates from unusual states (extreme glucose).
net_arch	[128,128,64]	Sufficient capacity for non-linear control. Not so large as to overfit. Bottleneck structure.
activation_fn	ReLU	Simple, effective, well-understood. No vanishing gradient issues.

Table 3: Complete hyperparameter summary with glucose control justifications

Part V

From Training to Inference

12 The Training Loop Explained

12.1 High-Level Training Flow

The training process in `train_ppo.py`:

Listing 11: Main training configuration and execution

```
1 def train(n_envs=4, total_steps=800_000, episode_days=1):
2     patient_id = 'adolescent#003'
3
4     env_kwargs = {
5         'patient_name': patient_id,
6         'reward_fun': paper_reward, # Or smart_reward
7         'seed': 42,
8         'episode_days': episode_days
9     }
10
11     # Create vectorized environments
12     train_env = make_vec_env(
13         CustomT1DEnv,
14         n_envs=n_envs,
15         env_kwargs=env_kwargs
16     )
17
18     # Create PPO model with all hyperparameters
19     model = PPO(
20         "MlpPolicy",
21         train_env,
22         verbose=1,
23         learning_rate=3e-4,
24         n_steps=2048,
25         batch_size=256,
26         n_epochs=10,
27         gamma=0.995,
28         gae_lambda=0.95,
29         clip_range=0.2,
30         ent_coef=0.01,
31         vf_coef=0.5,
32         max_grad_norm=0.5,
33         normalize_advantage=True,
34         policy_kwargs=dict(
35             net_arch=dict(pi=[128, 128, 64], vf=[128, 128, 64]),
36             activation_fn=torch.nn.ReLU
37         ),
38         tensorboard_log="ppo_smart_tensorboard",
39         device="auto"
40     )
```

```

41
42 # Train!
43 model.learn(
44     total_timesteps=total_steps,
45     callback=callback,
46     progress_bar=True
47 )

```

12.2 What Happens During `model.learn()`

Internally, Stable Baselines3 executes the PPO algorithm:

1. **Collect rollouts:** Each of 4 environments runs for 2048 steps
2. **Compute advantages:** Using GAE with $\gamma = 0.995$, $\lambda = 0.95$
3. **Normalize advantages:** Subtract mean, divide by std
4. **Shuffle data:** Randomize order of transitions
5. **Repeat 10 epochs:**
 - (a) Split into minibatches of 256
 - (b) For each minibatch:
 - i. Forward pass through policy network
 - ii. Compute probability ratios
 - iii. Compute clipped objective
 - iv. Compute value loss
 - v. Compute entropy bonus
 - vi. Backpropagate total loss
 - vii. Update parameters with Adam
 - viii. Clip gradients if necessary
6. **Update old policy:** $\theta_{old} \leftarrow \theta$
7. **Repeat** until `total_timesteps` reached

12.3 Monitoring Training: Callbacks

Listing 12: Training callbacks configuration

```

1 # Checkpoint callback: Save periodically
2 checkpoint_callback = CheckpointCallback(
3     save_freq=max(100_000 // n_envs, 1),
4     save_path="checkpoints_smart",
5     name_prefix=f"ppo_{patient_id}_smart"
6 )
7
8 # Evaluation callback: Track best model
9 eval_callback = EvalCallback(

```

```

10     eval_env,
11     best_model_save_path="checkpoints_smart",
12     log_path="eval_logs_smart",
13     eval_freq=max(50_000 // n_envs, 1),
14     n_eval_episodes=5,
15     deterministic=True
16 )

```

The EvalCallback is crucial:

- Every 50,000 steps, pause training
- Run 5 episodes with deterministic policy
- If mean reward improves, save as `best_model.zip`

This ensures we keep the best-performing checkpoint even if later training degrades.

13 Inference: How the Agent Makes Decisions

13.1 From Observation to Action

When deployed (in `main.py`), the trained agent operates deterministically:

Listing 13: Inference loop in `main.py`

```

1  # Load trained model
2  model = PP0.load(model_path)
3
4  # Get initial observation
5  obs, _ = env_rl.reset()
6
7  # Run simulation
8  for step in range(max_steps):
9      # THE KEY LINE: Get action from policy
10     action, _ = model.predict(obs, deterministic=True)
11
12     # Execute action in environment
13     obs, reward, terminated, truncated, info = env_rl.step(action)
14
15     # Record data
16     rl_bg_true.append(env_rl.env.patient.observation.Gsub)
17
18     if terminated or truncated:
19         break

```

13.2 Inside `model.predict()`

When `model.predict(obs, deterministic=True)` is called:

1. **Input:** Observation vector (27 floats)

```

obs = [0.7, 0.71, 0.69, ..., # 12 normalized glucose values
       0.1, 0.08, 0.12, ..., # 12 normalized insulin values
       0.4, 0.0, 0.0]       # 3 carb buckets

```

2. **Forward pass:** Through policy network

```

Layer 1: 27 -> 128 (ReLU)
Layer 2: 128 -> 128 (ReLU)
Layer 3: 128 -> 64 (ReLU)
Output:  64 -> 1 (linear)

```

3. **Output:** Mean of Gaussian policy

```

mu = 0.35 # Raw network output

```

4. **Deterministic mode:** Return μ directly (no sampling)

```

action = 0.35 # In range [-1, 1]

```

13.3 Action to Insulin: The Exponential Mapping

The raw action $a \in [-1, 1]$ must be converted to insulin rate:

Listing 14: Action mapping in custom_env.py step function

```

1 def step(self, action):
2     # 1. Clip action to valid range
3     a = np.clip(float(action[0]), -1.0, 1.0)
4
5     # 2. Exponential mapping to insulin rate
6     # I_max = 0.05 U/min, eta = 4.0
7     rate_u_min = self.I_max * np.exp(self.eta * (a - 1.0))
8
9     # 3. Create simglucose action
10    act = Action(basal=rate_u_min, bolus=0)
11
12    # 4. Step the underlying simulation
13    obs, reward, done, info = self.env.step(act, reward_fun=self.
        reward_fun)

```


Exponential Mapping Analysis

The mapping $I(a) = I_{max} \cdot e^{\eta(a-1)}$ with $I_{max} = 0.05$, $\eta = 4$:

Action a	Calculation	Insulin Rate (U/min)
-1	$0.05 \cdot e^{4(-2)} = 0.05 \cdot e^{-8}$	0.000017 (≈ 0)
0	$0.05 \cdot e^{4(-1)} = 0.05 \cdot e^{-4}$	0.00092
0.5	$0.05 \cdot e^{4(-0.5)} = 0.05 \cdot e^{-2}$	0.0068
1	$0.05 \cdot e^{4(0)} = 0.05 \cdot e^0$	0.05 (max)

Why Exponential Mapping?

1. **Safety:** $a = -1$ gives nearly zero insulin (basal suspension for hypoglycemia)
2. **Range:** $a = 1$ gives maximum safe rate
3. **Fine control:** Most actions are in low-dose region (where precision matters)
4. **Unbounded output from bounded input:** Natural for neural networks

14 Reward Engineering: Defining “Good” Control

14.1 The Paper Reward (Risk Index)

Listing 15: Risk index based reward function

```
1 def compute_risk_index(BG):
2     """Magni et al. Risk Index"""
3     bg_val = max(BG, 1.0) # Prevent log(0)
4     f_bg = 1.509 * ((np.log(bg_val) ** 1.084) - 5.381)
5     ri = 10 * (f_bg ** 2)
6     return ri
7
8 def paper_reward(bg_hist, **kwargs):
9     """Paper Eq. 6: Negative risk index"""
10    g_next = bg_hist[-1]
11    if g_next <= 39.0:
12        return -15.0 # Severe penalty for extreme hypoglycemia
13    raw_ri = compute_risk_index(g_next)
14    return -1.0 * (raw_ri / 100.0)
```

Risk Index Properties

The Magni et al. risk index is clinically validated:

- Symmetric risk for hypo and hyperglycemia at different glucose levels
- Minimum risk at euglycemic levels (≈ 100 -120 mg/dL)
- Asymptotically increasing risk at extremes

Approximate values:

Glucose (mg/dL)	Risk Index
50	22.7 (dangerous hypo)
70	5.8 (low)
100	0.5 (optimal)
150	2.1 (elevated)
200	6.3 (high)
300	17.2 (very high)

14.2 The Smart Reward (Enhanced)

Listing 16: Enhanced reward function with zones

```
1 def smart_reward(bg_hist, **kwargs):
2     """Enhanced reward with explicit zone bonuses"""
3     bg = bg_hist[-1]
4
5     # Extreme hypoglycemia: Severe penalty
6     if bg <= 40:
7         return -100.0
8
9     # Target range (tight): Bonus reward
10    if 70 <= bg <= 150:
11        return 1.0
12
13    # Acceptable range: Small bonus
14    if 150 < bg <= 180:
15        return 0.5
16
17    # Outside target: Scaled risk penalty
18    risk = compute_risk_index(bg)
19    multiplier = 2.0 if bg > 250 else 1.0 # Extra penalty for very
20    high
    return -1.0 * (risk / 50.0) * multiplier
```

Why Two Reward Functions?

Paper reward (risk index):

- Clinically validated, used in research
- Smooth gradient everywhere
- But: Agent might not “try hard” to stay in target if risk is low outside

Smart reward (zone-based):

- Explicit positive reward for target range
- Encourages active maintenance, not just risk avoidance
- Faster learning in practice
- But: Less theoretically grounded

15 Robustness Testing: Stress-Testing the Agent

The `test_robustness.py` script evaluates the trained agent on challenging scenarios:

Listing 17: Stress test scenarios

```
1 def get_scenarios():
2     """Define stress test scenarios"""
3     scenarios = {}
4
5     # 1. High Carbs: Can the agent handle large meals?
6     high_carb_meals = [
7         (timedelta(hours=8), 80),    # Large breakfast
8         (timedelta(hours=13), 110), # Huge lunch
9         (timedelta(hours=20), 90)    # Heavy dinner
10    ]
11    scenarios['High_Carbs'] = CustomScenario(...)
12
13    # 2. Missed Lunch: Does the agent avoid hypoglycemia?
14    missed_lunch_meals = [
15        (timedelta(hours=8), 40),
16        # No lunch!
17        (timedelta(hours=20), 60)
18    ]
19    scenarios['Missed_Lunch'] = CustomScenario(...)
20
21    # 3. Late Dinner: Overnight hypoglycemia risk?
22    late_dinner_meals = [
23        (timedelta(hours=8), 40),
24        (timedelta(hours=13), 60),
25        (timedelta(hours=22, minutes=30), 70) # Very late
26    ]
27    scenarios['Late_Dinner'] = CustomScenario(...)
```

```

28
29     # 4. Random Chaos: General robustness
30     # 4 random meals throughout the day
31     scenarios['Random_Chaos'] = CustomScenario(...)
32
33     return scenarios

```

What Each Scenario Tests

1. **High Carbs:** Tests maximum insulin delivery capacity and post-prandial control
2. **Missed Lunch:** Tests if agent reduces insulin when expected carbs don't arrive (hypoglycemia prevention)
3. **Late Dinner:** Tests overnight insulin adjustment (dawn phenomenon, nocturnal hypoglycemia)
4. **Random Chaos:** Tests generalization beyond training distribution

Part VI

Conclusion and Future Directions

16 Summary of Key Insights

16.1 Why PPO Works for Glucose Control

1. **Continuous control:** Natural handling of insulin as continuous variable
2. **Safety through clipping:** Prevents dangerous policy changes
3. **Sample efficiency:** Multiple epochs reuse expensive simulation data
4. **Stability:** Monotonic improvement guarantees
5. **Proactivity:** Meal announcements enable pre-bolusing

16.2 Critical Implementation Details

1. **State augmentation:** History buffers restore Markov property
2. **Normalization:** Essential for stable neural network training
3. **Exponential action mapping:** Safe, bounded insulin delivery
4. **Reward shaping:** Zone-based rewards accelerate learning
5. **Episode termination:** Catastrophic glucose values end episodes with penalty

17 Limitations and Future Work

17.1 Current Limitations

- Single patient training (not generalizing across patients)
- Perfect meal announcements (real patients don't always pre-announce)
- No sensor noise during training (CGM in simulation is idealized)
- Fixed meal patterns during training

17.2 Future Directions

- **Multi-patient training:** Train on diverse patient population
- **Meal uncertainty:** Random or absent meal announcements
- **Sensor noise injection:** More realistic CGM simulation
- **Transfer learning:** Pre-train on simulation, fine-tune on real data
- **Safety constraints:** Constrained RL to enforce hard bounds

18 Final Remarks

This report has provided a comprehensive technical analysis of PPO for automated insulin delivery. The algorithm's combination of sample efficiency, stability, and safety makes it particularly well-suited for medical control applications. The key insight—that clipping the objective function can replace complex KL constraints while maintaining policy stability—enables practical deployment of RL in safety-critical systems.

The hyperparameters chosen reflect deep understanding of both the algorithm mechanics and the clinical domain. The high discount factor accounts for insulin's long duration of action. The conservative clip range prioritizes safety over aggressive learning. The entropy coefficient balances exploration against the risks of random actions in a medical context.

Through careful state design, reward engineering, and hyperparameter tuning, we have demonstrated that a PPO agent can learn effective glucose control policies that outperform traditional methods on challenging virtual patients. This represents a significant step toward fully automated artificial pancreas systems.