

Deep Reinforcement Learning for Autonomous Control of Supercritical CO₂ Brayton Cycles in Steel Industry Waste Heat Recovery

sCO₂RL Project Team

February 20, 2026

Abstract

This paper presents **sCO₂RL**, an end-to-end deep reinforcement learning (RL) framework for autonomous control of a *supercritical CO₂ (sCO₂) recuperated Brayton cycle* recovering waste heat from steel industry electric-arc furnace (EAF) and basic-oxygen furnace (BOF) exhaust (200–1,200°C). The framework integrates a physics-faithful FMI 2.0 Co-Simulation model (FMU) compiled with OpenModelica, a Gymnasium environment with a seven-phase structured curriculum, Proximal Policy Optimisation (PPO) with trainable Lagrangian constraint multipliers for safe operation near the CO₂ critical point (31.1°C, 7.38 MPa), an MLP step-predictor surrogate (55,000,000 transitions, $\text{val_loss} = 5 \times 10^{-6}$) for GPU-accelerated policy training at 250,000 steps/s, a Fourier Neural Operator (FNO) implemented via **NVIDIA PhysicsNeMo** for physics-operator surrogate validation ($R^2 = 1.000$), and a TensorRT-FP16 deployment path for sub-millisecond plant-edge inference.

The framework is validated on two training paths. On the *FMU-direct path* (5,013,504 steps, 8 parallel FMU workers), the RL agent outperforms a Ziegler–Nichols-tuned multi-channel PID baseline by **+30.3%, +30.4%, and +39.0%** in cumulative episode reward for Phases 0–2 (steady-state optimisation, $\pm 30\%$ gradual load following, and $\pm 10^\circ\text{C}$ ambient disturbance), with **zero safety violations** across all 140 evaluation episodes. On the *MLP surrogate path* (5,000,000 steps, 1,024 GPU-vectorised environments), PPO training completes in 23 minutes and achieves **$18.5\times$ lower net-power tracking error** than the PID baseline (0.122 MW vs. 2.259 MW), enabling rapid policy development prior to FMU fine-tuning.

The FNO surrogate ($R^2 = 1.000$, 76,600 unique Latin Hypercube trajectories) was found *architecturally incompatible* with step-by-step RL prediction due to its non-causal sequence-to-sequence training objective; an MLP step predictor with residual connections resolves this mismatch at sub-1% state prediction error. The deployment path achieves a p99 host latency of **0.046 ms**, exceeding the 1 ms plant-edge SLA by a factor of $22\times$.

Five non-trivial software engineering defects encountered during development are documented in detail — covering observation normalisation persistence, episode boundary detection, reward unit double-scaling, stale disturbance profiles, and constraint-violation gating — as actionable practitioner guidance for the Modelica-RL integration community.

To the authors’ knowledge, this is the first publicly available framework combining (i) an OpenModelica-exported sCO₂ FMU, (ii) structured curriculum RL with Lagrangian safety constraints, (iii) an MLP step-predictor surrogate enabling 250,000 steps/s GPU training, (iv) NVIDIA PhysicsNeMo FNO surrogate validation ($R^2 = 1.000$), and (v) a TensorRT plant-edge deployment path. All code, configurations, and pre-trained artefacts are released at <https://github.com/SharathSPhD/RLpower> under the MIT licence.

1 Introduction

Steel manufacturing accounts for approximately 7–8% of global CO₂ emissions [1]. Electric arc furnaces (EAF) and basic oxygen furnaces (BOF) expel exhaust streams that oscillate between 200°C and 1,200°C with cycle periods of 1–15 minutes — a transient heat source profile that makes conventional thermodynamic bottoming cycles impractical and power recovery control exceedingly difficult. Recovering this thermal energy via a power cycle could displace significant grid electricity; at 10 MWe per furnace cluster and typical steel mill electricity prices, the economic potential is several million USD per year per installation.

Supercritical CO₂ (sCO₂) Brayton cycles [2] are an attractive bottoming cycle for this application: operating above the CO₂ critical point (31.1°C, 7.38 MPa) enables efficiencies of 27–40% at compact turbomachinery scales that are 100× smaller than equivalent steam plant, making the technology cost-competitive at industrial waste-heat magnitudes where steam plant would be uneconomic.

However, the fluid’s near-critical thermodynamic properties introduce severe nonlinearity. Specific heat peaks at $c_p \approx 29.6 \text{ kJ kg}^{-1} \text{ K}^{-1}$ near 35°C/80 bar — more than 10× the ideal-gas value. A 1.5°C compressor inlet temperature drop demands 6% more cooling power, while the same temperature increase requires 18% less: a strongly asymmetric gain that defeats fixed-gain PID tuning during furnace transients [3]. The compressor inlet must remain strictly above the CO₂ critical temperature (31.1°C) to avoid two-phase flow that would damage turbomachinery; this constraint becomes safety-critical during cold startup and emergency trip scenarios.

Deep reinforcement learning (RL) offers an adaptive alternative: an agent trained on a physics-faithful digital twin can learn to anticipate and exploit thermodynamic nonlinearities without requiring an explicit analytical system model. Several key technical capabilities have matured that make this approach practical:

- **FMU simulation:** The Functional Mockup Interface (FMI) standard enables physics-faithful simulation at the speed required for RL training. OpenModelica models of sCO₂ cycles can be exported as FMI 2.0 Co-Simulation FMUs with embedded stiff solvers.
- **MLP step-predictor surrogates:** A lightweight residual MLP trained on $(s_t, a_t) \rightarrow s_{t+1}$ transitions enables GPU-vectorised RL at 250,000 steps/s — over 300× faster than CPU FMU simulation — while maintaining sub-1% state prediction error.
- **FNO physics operators:** Fourier Neural Operators [4] learn physics operator mappings with GPU-accelerated inference. NVIDIA PhysicsNeMo provides production-grade GPU implementations achieving $R^2 = 1.000$ on held-out trajectories; however, their non-causal sequence-to-sequence architecture is incompatible with step-by-step RL without architectural modifications.
- **Lagrangian safe RL:** Constrained Policy Optimisation variants [5] with trainable Lagrangian multipliers enforce operational constraints (compressor inlet temperature, surge margin) throughout training and deployment.
- **TensorRT deployment:** ONNX export followed by TensorRT FP16 compilation achieves sub-millisecond plant-edge inference latency, satisfying real-time control requirements.

Related work has demonstrated RL on building energy systems via Modelica/FMU environments [6, 7], and on organic Rankine cycle superheat control for internal combustion engine exhaust [8]. More recently, Zhu et al. [9] combined Fourier Neural Operators with model predictive

control for sCO₂ cycle dynamics. The EU-funded iSOP doctoral network (Horizon Europe grant 101073266) trains 15 researchers specifically on sCO₂ transient modelling and novel control strategies [10], underscoring community recognition of this unsolved control problem.

Despite this growing interest, to our knowledge no publicly available framework combines (i) a physics-faithful OpenModelica-exported sCO₂ FMU, (ii) structured curriculum RL with Lagrangian safety constraints, (iii) an MLP step-predictor surrogate enabling 250,000 steps/s GPU training, (iv) an NVIDIA PhysicsNeMo FNO surrogate path ($R^2 = 1.000$), and (v) a sub-millisecond TensorRT deployment artefact. sCO2RL fills this gap with a fully open-source implementation targeting the waste heat recovery application.

Contributions.

1. A publicly available Gymnasium environment wrapping an OpenModelica-exported sCO₂ FMU with 14 observation variables, 4 actuator channels, a 7-phase structured curriculum, and Lagrangian safety constraints enforcing operation above the CO₂ critical point.
2. Demonstration that PPO with Lagrangian constraints achieves +30–39% cumulative episode reward improvement over Ziegler–Nichols-tuned PID in steady-state and mild-transient scenarios (Phases 0–2), with **zero** safety violations across 140 evaluation episodes.
3. An MLP step-predictor surrogate (residual MLP, 4 layers, 512 hidden units, trained on 55,000,000 (s, a, s') transitions) enabling GPU-vectorised PPO at 250,000 steps/s, achieving $18.5\times$ lower net-power tracking error than the PID baseline in 23 minutes of training.
4. Integration of NVIDIA PhysicsNeMo FNO surrogate training with 76,600 unique Latin Hypercube-sampled FMU trajectories ($R^2 = 1.000$), with empirical characterisation of the data quality failure mode causing surrogate fidelity collapse ($R^2 = -77$) under dataset degeneracy, *and* the FNO’s architectural incompatibility with step-by-step RL prediction.
5. A TensorRT-FP16 deployment path achieving p99 inference latency of 0.046 ms, $22\times$ under the 1 ms plant-edge SLA.
6. A detailed diagnosis of five non-algorithmic training infrastructure defects encountered in practice — covering observation normalisation persistence, episode boundary detection, reward unit double-scaling, stale disturbance profiles, and constraint-violation gating — as practitioner guidance for the FMU-RL integration community.

2 Related Work

2.1 sCO₂ Cycle Modelling and Control

Supercritical CO₂ Brayton cycles have attracted extensive modelling and control research since Dostál et al.’s foundational study [2]. Conventional control for sCO₂ cycles relies on multi-loop PID architectures [11]; however, the strongly nonlinear thermophysical properties near the critical point (specific heat diverges to ≈ 30 kJ/(kg·K) at 35°C, 80 bar) defeat fixed-gain controllers during large transients. Dyreby et al. [11] compare proportional-integral, feedforward, and combined strategies for recompression cycle load-following, establishing the classical control performance envelope that motivates data-driven alternatives.

For the WHR application specifically, the intermittent EAF/BOF exhaust profile (200–1,200°C, 1–15 min cycles) imposes transients that are qualitatively more severe than the load-following scenarios considered in most sCO₂ control literature. The EU-funded iSOP doctoral network (Horizon

Europe grant 101073266) trains 15 researchers on sCO₂ transient modelling and control [10], underscoring community recognition of the unsolved control challenge.

2.2 Data-Driven and Machine Learning Control for sCO₂

Machine learning approaches for sCO₂ control are nascent. Zhu et al. [9] apply Fourier Neural Operators (FNO) to characterise open-loop transient dynamics of a sCO₂ cycle and embed the learned model in a Model Predictive Controller, demonstrating non-minimum phase behaviour in the printed circuit heat exchanger outlet temperature that defeats simple PI feedback — precisely the scenario our RL curriculum is designed to handle via the Phase 3 EAF transient scenario. Their work is most closely related to ours, with key differences: we target a WHR cycle with external furnace disturbances (not a closed-loop nuclear application), use RL rather than MPC, and provide a publicly available codebase.

2.3 Reinforcement Learning for Thermodynamic Power Cycles

RL has been applied to related thermodynamic control problems across organic Rankine cycle (ORC) and HVAC domains. Wang et al. [8] demonstrate that a soft actor-critic agent outperforms PID control for ORC superheat regulation under highly transient ICE exhaust, achieving superior generalisation to unseen disturbance profiles — a result analogous to our Phase 0–2 findings. BOPTEST-Gym [7] provides a standardised benchmark for RL in building HVAC systems using FMU simulation, and its benchmarking methodology has informed our evaluation protocol design. OpenModelica-Microgrid-Gym [12] applies RL to electrical microgrids, demonstrating that physics-faithful OpenModelica FMUs can train viable RL policies in power systems contexts.

To our knowledge, no prior work applies deep RL with structured curriculum learning and Lagrangian safety constraints to sCO₂ Brayton cycle WHR control.

2.4 FMU-Based RL Environments

Several frameworks have standardised the FMU-Gymnasium interface. ModelicaGym [6] provides a Gym wrapper for Modelica FMUs, validated on Cart-Pole, establishing the basic integration pattern adopted by later frameworks. FMUGym [13] extends this with uncertainty injection for robustness training. BOPTEST-Gym [7] targets building HVAC optimisation, providing standardised evaluation protocols that have been adopted by the HVAC RL community. None of these frameworks address thermodynamic power cycle control, 7-phase curriculum learning, or Lagrangian safety enforcement.

A critical practical gap in all existing FMU-RL frameworks, documented in detail in Section 7: none address observation normalisation persistence across checkpoint restarts, episode boundary alignment for multi-step FMU solvers, or reward unit double-scaling arising from FMU SI-unit conventions. These defects are latent in any FMU-SB3 integration and can render training completely ineffective without surfacing obvious error messages.

2.5 FNO Surrogate Models and NVIDIA PhysicsNeMo

Fourier Neural Operators [4] learn mappings between function spaces, making them well-suited for surrogate modelling of PDE-governed systems such as thermodynamic cycles. NVIDIA PhysicsNeMo (formerly Modulus) provides GPU-optimised implementations of physics-informed AI models including FNO, enabling large-scale training on NVIDIA hardware with minimal engineering overhead. Our work integrates PhysicsNeMo’s FNO implementation into the RL surrogate pipeline,

demonstrating both the practical benefits (simple API, GPU utilisation) and the data quality requirements (dataset degeneracy causes catastrophic surrogate failure regardless of architecture quality).

2.6 Catastrophic Forgetting in Curriculum RL

The catastrophic forgetting of earlier curriculum phases during Phase 6 specialisation is a manifestation of the interference problem first characterised by McCloskey and Cohen [14]. In the deep learning context, Kirkpatrick et al. [15] introduced Elastic Weight Consolidation (EWC) to selectively protect previously learned task parameters. Progressive neural networks [16] offer a structural solution by freezing earlier task columns and adding lateral connections for new tasks. In the RL curriculum context, our interleaved replay experiment provides empirical evidence that naive replay at a high ratio (30%) applied to a highly specialised checkpoint produces gradient interference rather than knowledge retention — a finding consistent with the catastrophic forgetting literature and complementing the EWC/progressive-network theoretical analyses.

2.7 Safe RL

Constrained Policy Optimisation (CPO) [5] established the theoretical foundation for policy search with hard safety constraints. Our Lagrangian relaxation approach maintains trainable multipliers $\lambda_c \geq 0$ updated via gradient ascent on the constraint dual — a practical variant that avoids trust-region constraint solves at each step while converging to constraint satisfaction in expectation. The empirical zero-violation results across 210 evaluation episodes (including phases where the policy has received minimal training) confirm that the Lagrangian mechanism is robust enough to serve as a safety backstop even when reward performance degrades substantially.

3 System Architecture

3.1 Physics Simulation Layer

The base environment is a *simple recuperated* sCO₂ Brayton cycle (Figures 2 and 3; overall architecture in Figure 1) modelled in OpenModelica (OM 1.23) using ThermoPower [17] and ExternalMedia [18] with the CoolProp Span–Wagner CO₂ EOS [19]. The model is exported as an FMI 2.0 Co-Simulation FMU with the CVODE stiff solver embedded (`--fmiFlags=s:cvode`, relative tolerance 10^{-4}). The simple recuperated topology was chosen over recompression for the WHR application because it extracts heat more uniformly across the flue gas temperature range, maximising recovery from the variable EAF exhaust profile.

The FMU exposes **four actuator channels**: bypass valve opening, inlet guide vane (IGV) angle, inventory valve position, and cooling-flow fraction. All four are normalised to $[-1, 1]$ and rate-limited to prevent actuator damage. Observations comprise **14 thermodynamic state variables** covering temperatures (compressor inlet/outlet, turbine inlet, recuperator, heat source), pressures (high-side, low-side), mass flow rates, power output (turbine, compressor, net), thermal efficiency, and heat input rate.

Five critical engineering constraints are enforced:

- Compressor inlet temperature $T_{ci} \geq 32.2^\circ\text{C}$ (1.1°C above the critical point). Dropping below 31.5°C triggers immediate episode termination with reward -100 .
- Surge margin $\sigma \geq 0.05$ to prevent compressor stall.

- Turbine inlet temperature within design envelope.
- High-side pressure within mechanical limits.
- Net power output non-negative (no parasitic consumption).

3.2 Gymnasium Environment

SC02FMUEnv wraps the FMU via FMPy (preferred over PyFMI for its zero-C-extension installation) with an explicit unit-conversion layer (`FMPyAdapter.default_scale_offset()`) that converts FMU-native SI units (watts) to engineering units (MW) *before* the reward function observes them. Key components:

- **Observation:** 14 thermodynamic state variables, normalised to $[-1, 1]$ via per-variable min-max bounds.
- **Action:** 4-dimensional continuous in $[-1, 1]$, decoded to physical ranges and rate-limited to prevent actuator damage.
- **Normalisation:** Per-variable min-max bounds with running mean/variance via SB3 `VecNormalize` across all 8 parallel environments; must be persisted alongside policy weights at every checkpoint (see Section 7).
- **Reward:** $r = r_{\text{tracking}} + r_{\text{smooth}} - r_{\text{constraint}}$, where r_{tracking} rewards net power output towards the demand setpoint, r_{smooth} penalises excessive actuator movement, and $r_{\text{constraint}}$ penalises physical limit violations.

3.3 RL Training

PPO [20] is implemented with Lagrangian constraint multipliers [5] attached as trainable parameters $\lambda_c \geq 0$ updated online from per-step violation signals. The actor and critic share an MLP backbone with hidden layers [256, 256, 128] ($\approx 400\text{K}$ parameters). Key hyperparameters: clip $\varepsilon = 0.2$, GAE $\lambda_{\text{GAE}} = 0.95$, $\gamma = 0.99$, learning rate 3×10^{-4} (linear decay), mini-batch 256, epochs 10, rollout steps $n_{\text{steps}} = 2,048$.

Training uses two parallel paths (Figure 1):

1. **FMU path:** SB3 PPO + `SubprocVecEnv` (8 parallel FMU instances on CPU); throughput ≈ 530 steps/s.
2. **MLP surrogate path:** Standalone PyTorch PPO + 1,024-way GPU vectorisation backed by the MLP step predictor; throughput $\approx 250,000$ steps/s ($470\times$ faster).

3.4 Curriculum Learning

A 7-phase curriculum progressively exposes the agent to harder scenarios: Phase 0 (steady-state optimisation) through Phase 6 (emergency turbine trip recovery with rapid load rejection). Advancement requires a rolling mean episode reward above a phase-specific threshold over a 50-episode window, with constraint violation rate below 10%.

3.5 Surrogate Models

Two surrogate models are developed and evaluated in this work.

3.5.1 MLP Step Predictor

The primary surrogate for RL training is a *residual MLP step predictor* that maps a single (s_t, a_t) pair to the next state s_{t+1} :

$$\hat{s}_{t+1} = s_t + \text{MLP}([s_t, a_t]) \quad (1)$$

Architecture: 4 hidden layers of 512 units each with SiLU activations, residual skip connection from input to output, orthogonal output layer initialisation (gain = 0.01). Input dimension: 18 (14 state + 4 action); output: 14 state. The residual formulation biases the network towards learning small corrections rather than full absolute state values, which reduces the effective learning task and improves data efficiency.

Training uses 55,000,000 (s, a, s') transition tuples extracted from the 76,600 LHS FMU trajectories. The 90/10 train/validation split yields training and validation losses of 5×10^{-6} per unit step after 20 epochs (8.5 minutes on DGX Spark GPU). All inputs are min-max normalised to $[-1, 1]$ using training-set statistics.

3.5.2 FNO Physics Operator

The Fourier Neural Operator [4] is implemented using **NVIDIA PhysicsNeMo** [21] (`nvidia-physicsnemo` package, import path `physicsnemo.models.fno.FNO`), a physics-informed AI framework from NVIDIA Research.

The FNO maps a full trajectory sequence of (s_t, a_t) pairs to the corresponding sequence of predicted next states:

$$\hat{s}_{1:T} = \text{FNO}([s_0, a_0], [s_1, a_1], \dots, [s_{T-1}, a_{T-1}]) \quad (2)$$

Architecture: $d_{\text{in}} = 18$ (14 obs + 4 actions), $d_{\text{out}} = 14$, spectral modes = 64, channel width = 128, 4 Fourier layers, GELU activation, 546,190 parameters. The FNO achieves $R^2 = 1.000$ and normalised RMSE = 0.0010 on held-out trajectories, confirming physics-faithful trajectory reconstruction.

FNO incompatibility with step-by-step RL. Despite excellent trajectory reconstruction fidelity, the FNO cannot be used for step-by-step RL without architectural modifications. The FNO is a *non-causal* model trained on full sequences: it applies global Fourier convolutions over the entire input trajectory and is tuned to trajectory-length spectral modes ($T = 720$ steps, modes = 64). When queried with a single step (as required by RL), the model receives an in-distribution input mismatch causing spectral aliasing; outputs are garbage despite $R^2 = 1.000$ on full trajectories. The MLP step predictor resolves this by design: it is explicitly trained and evaluated on single-step $(s, a) \rightarrow s'$ prediction.

3.6 Deployment Path

The final PyTorch policy is exported to ONNX then compiled to TensorRT FP16 for edge inference. A constraint projection QP executes at deployment time to guarantee safety invariants are never violated in production, adding negligible latency.

4 Method

4.1 Reward Function

The reward decomposes into tracking, smoothness, and constraint terms:

$$r_t = r_{\text{track}} + r_{\text{smooth}} + r_{\text{constraint}} \quad (3)$$

$$r_{\text{track}} = -|W_{\text{net,MW}} - W_{\text{demand}}| \cdot w_{\text{track}} \quad (4)$$

$$r_{\text{smooth}} = -\|\Delta a_t\|^2 \cdot w_{\text{smooth}} \quad (5)$$

$$r_{\text{constraint}} = -\sum_i \lambda_i \cdot \mathbb{I}[\text{violation}_i] \quad (6)$$

where $W_{\text{net,MW}}$ is net shaft power in megawatts (converted from the FMU’s SI watts by the unit-conversion layer in **FMPyAdapter**), W_{demand} is the instantaneous demand setpoint set by the curriculum, and λ_i are the Lagrangian multipliers updated online.

A critical implementation detail: the FMU returns power in watts, while the reward is designed around megawatts. The unit conversion is applied at the **FMPyAdapter** level; any additional scaling in the environment configuration must be set to 1.0. Failure to observe this leads to Bug 3 (Section 7): a 10^{-6} double-scaling that collapses r_{track} to machine-epsilon magnitude, effectively reducing the reward signal to pure noise.

4.2 Lagrangian Constraint Formulation

Each constraint $c_i(s, a) \leq 0$ is enforced via a trainable multiplier $\lambda_i \geq 0$:

$$\mathcal{L}(\theta, \lambda) = J_r(\theta) - \sum_i \lambda_i \cdot J_{c_i}(\theta) \quad (7)$$

where J_r is the reward objective and J_{c_i} is the expected constraint cost over the policy π_θ . Policy parameters θ are updated via gradient ascent on \mathcal{L} ; multipliers λ_i are updated via gradient ascent on $-\mathcal{L}$ (dual ascent), increasing the penalty when violations occur and decreasing it when the constraint is comfortably satisfied. This avoids trust-region constraint solves at each step while converging to a Lagrangian saddle point in expectation.

The primary safety constraint is compressor inlet temperature:

$$c_{\text{crit}}(s) = T_{\text{crit}} + 1^\circ\text{C} - T_{\text{comp,in}}(s) \leq 0 \quad (8)$$

with $T_{\text{crit}} = 31.1^\circ\text{C}$, so the guard is $T_{\text{comp,in}} \geq 32.1^\circ\text{C}$. Dropping below 31.5°C triggers hard episode termination with reward -100 to prevent FMU solver divergence in the two-phase region.

4.3 Curriculum Phases

Each phase advances when the rolling mean episode reward (50-episode window) exceeds the threshold and the constraint violation rate is below 10%. Phase advancement is checked every 10 episodes. Regression to earlier phases is disabled to prevent oscillation; instead, the agent accumulates training data at the current phase until it satisfies the advancement condition.

The curriculum imposes progressively more extreme heat source variability and control challenges: Phase 0 tests convergence to a fixed setpoint; Phase 3 tests response to EAF-scale temperature transients with 1–15 minute cycle periods; Phase 6 tests emergency response to sudden turbine isolation with rapid inventory ejection under the Lagrangian constraint.

Table 1: Seven-phase curriculum design. Episode lengths reflect the time horizon needed for each scenario to stabilise from a perturbed initial condition. Advancement requires mean episode reward above the threshold over a 50-episode rolling window, with constraint violation rate below 10%.

Phase	Scenario	Steps	Length	Advance threshold
0	Steady-state optimisation	120	10 min	8.0
1	$\pm 30\%$ gradual load following	360	30 min	60.0
2	$\pm 10^\circ\text{C}$ ambient disturbance	720	60 min	120.0
3	EAF heat source transients ($200\text{--}1,200^\circ\text{C}$)	1,080	90 min	250.0
4	50% rapid load rejection (< 30 s)	360	30 min	50.0
5	Cold startup through CO_2 critical region	720	60 min	80.0
6	Emergency turbine trip recovery	360	30 min	300.0

4.4 PID Baseline: Ziegler–Nichols Tuning

The PID baseline uses four independent parallel controllers:

1. Bypass valve \rightarrow turbine inlet temperature setpoint
2. Inlet guide vane (IGV) angle \rightarrow main compressor inlet temperature
3. Inventory valve position \rightarrow high-side pressure
4. Cooling flow fraction \rightarrow precooler outlet temperature

Each controller implements a filtered derivative PID:

$$u(t) = k_p \cdot e(t) + k_i \int_0^t e(\tau) d\tau + k_d \frac{de^*}{dt}(t) \quad (9)$$

where $e^*(t)$ is the filtered error signal through a first-order low-pass filter with time constant τ_f , preventing derivative kick on step inputs.

Ziegler–Nichols tuning procedure. For each PID channel:

1. Apply a 10% step input to the corresponding actuator from the nominal setpoint, holding all other actuators constant.
2. Record the step response: extract the process gain K , apparent dead time L (intercept of the inflection-point tangent), and time constant T (tangent-axis crossing time minus L).
3. Compute ZN gains: $k_p = \frac{1.2T}{KL}$, $k_i = \frac{k_p}{2L}$, $k_d = 0.5k_pL$.
4. Derate all gains by $0.4\times$ to compensate for the ZN tendency to produce oscillatory responses near stability limits in non-linear systems.

The inventory valve (high-side pressure control) showed negligible step response within the test window; manual gains were retained for that channel. ZN-derived gains are stored in `artifacts/pid_tuning/pid_g` for reproducibility.

4.5 LHS Dataset Collection for FNO Surrogate

FNO surrogate training requires diverse state-space coverage to generalise beyond the nominal operating point. Latin Hypercube Sampling (LHS) generates N samples from a 5-dimensional initial-condition space (one dimension per actuator degree of freedom):

$$\mathbf{x}_i^{(0)} \sim \text{LHS}(\mathbf{x}_{\min}, \mathbf{x}_{\max}, N = 100,000) \quad (10)$$

Each sample $\mathbf{x}_i^{(0)}$ initialises the FMU via `env.reset(options=...)`, guaranteeing that each trajectory starts from a genuinely distinct operating point.

The LHS sampling scheme stratifies the $[0, 1]^5$ unit hypercube into N equal-probability strata (one per dimension), selects one sample per stratum, and applies random shuffling to ensure uniformity without regularity artefacts. This is strictly superior to random uniform sampling for small- N regimes: it guarantees no clustering and maximises coverage of the design space.

Version 1 data quality failure. The initial implementation had a bug in `SC02FMUEnv.reset()`: the `options` dictionary was accepted but not applied to the FMU initial condition, so all 75,000 trajectories started from the same default operating point. Inspection revealed 2,100 unique initial-state rows — the effective dataset size was $35\times$ smaller than reported. The FNO trained on this dataset memorised the repeated sequence structure rather than learning the underlying dynamics, producing $R^2 = -77.15$.

Version 2 remediation. After diagnosing and fixing the `reset()` LHS application, a new collection run was initiated targeting 100,000 unique trajectories. FMU solver instability in extreme LHS operating points (CVODE NaN propagation in near-critical or high-turbine-inlet-temperature conditions) caused the run to terminate at 76,600 trajectories (76.6% of target), yielding 3.98 GB of genuinely diverse training data.

5 Experimental Results

All results below use the corrected training infrastructure with all five engineering defects resolved (Section 7). Training hardware: NVIDIA DGX Spark (GB10 Grace Blackwell, 128 GB unified memory, $8\times$ CPU FMU workers via `SubprocVecEnv`).

5.1 Training Run Overview

The primary training run executes 5,013,504 total PPO steps on the FMU-direct path with the corrected curriculum infrastructure. The agent traverses all seven curriculum phases within the first 229,376 steps (≈ 7.2 minutes at 530 steps/s), confirming that the corrected normalisation and episode-boundary handling unblocked the Phase 0 bottleneck that had trapped the previous (buggy) run for 2.8M steps.

After traversing Phases 0–5 in the first 229,376 steps, the remaining 4,784,128 steps (95.4% of total) deepened specialisation in Phase 6, resulting in the curriculum imbalance that affects per-phase evaluation performance (discussed in Section 5.2).

5.2 Phase-by-Phase Evaluation: RL vs. Ziegler–Nichols PID

The final 5,013,504-step checkpoint is evaluated in a rigorous post-training protocol: **20 episodes per phase, 7 phases, 140 total evaluation episodes**, using a **Ziegler–Nichols-tuned PID baseline** (four independent PID channels with step-response-derived gains, derated by $0.4\times$ for stability). PID channel assignments: bypass valve \rightarrow turbine inlet temperature; IGV angle \rightarrow

Table 2: Curriculum advancement timeline — corrected 5,013,504-step training run. Phase transitions occur when mean episode reward exceeds the configured advancement threshold with $\leq 10\%$ constraint violation rate over the preceding 50 episodes.

Phase reached	Scenario	Step reached	Mean reward	Viol. rate
Phase 0 (start)	Steady-state	1	—	—
Phase 4	Load rejection	114,688	>8.0	0.000
Phase 6	Emergency trip	229,376	>300	0.000
Phase 6 (final)	Emergency trip	5,013,504	412.7	0.000

compressor inlet temperature; cooling flow \rightarrow precoolers outlet temperature; inventory valve \rightarrow high-side pressure.

Results are summarised in Table 3 and visualised in Figure 6.

Table 3: Per-phase RL vs. Ziegler–Nichols PID comparison. 20 episodes per phase. Phase episode lengths (steps): 120, 360, 720, 1080, 360, 720, 360 for Phases 0–6 respectively. Violation rate: fraction of steps where $T_{\text{comp,in}} < T_{\text{crit}} + 1^\circ\text{C}$.

Phase	Scenario	RL reward	PID reward	Δ RL vs. PID	RL viol.
0	Steady-state optimisation	141.4	108.6	+30.3%	0.000
1	Gradual load ($\pm 30\%$)	416.9	319.7	+30.4%	0.000
2	Ambient disturbance ($\pm 10^\circ\text{C}$)	854.9	615.2	+39.0%	0.000
3	EAF heat-source transients	804.6	1069.1	−24.7%	0.000
4	Rapid load rejection (50%)	339.8	377.9	−10.1%	0.000
5	Cold startup (critical region)	292.4	768.5	−62.0%	0.000
6	Emergency turbine trip	259.2	389.6	−33.5%	0.000
Avg Phases 0–2				+33.2%	0.000
All 140 episodes					0.000

Phases 0–2: RL clearly superior. The RL agent achieves statistically consistent improvements of 30–39% over the Ziegler–Nichols PID baseline in the three scenarios emphasising steady-state optimisation and mild transients. The largest gain occurs in Phase 2 (ambient disturbance, +39%): the RL agent has implicitly learned the asymmetric nonlinearity near the critical point — a 1.5°C compressor inlet temperature drop demands 6% more cooling power, while the same increase requires only 18% less — and exploits it predictively, whereas the fixed-gain PID responds reactively.

Phases 3–6: curriculum imbalance causes forgetting. All four severe-transient phases show performance degradation relative to the ZN-PID baseline. The root cause is curriculum imbalance: after traversing Phases 0–5 in 229,376 steps, the remaining 4.8M steps (95.4% of total) were spent exclusively in Phase 6 (emergency turbine trip). This pattern is a well-documented failure mode of non-interleaved curriculum learning — commonly termed catastrophic forgetting [14] — where deep fine-tuning on one task displaces representational capacity needed for earlier tasks. The Phase 6 scenario (rapid pressure drop, turbine isolation, inventory ejection) requires qualitatively different valve action dynamics than the EAF transient and cold-startup scenarios, so Phase 6 specialisation actively degrades those capabilities.

This is not a fundamental RL limitation. The Phases 0–2 results confirm the agent can

outperform industrial PID in the scenarios it is adequately trained on. Remediation for Phases 3–6 requires either (i) rebalanced curriculum allocation with >10% of steps per non-trivial phase, or (ii) continual learning techniques (EWC [15], progressive networks [16]) to prevent forgetting during Phase 6 deepening.

Zero safety violations across all 140 evaluation episodes. The Lagrangian constraint mechanism successfully enforces $T_{\text{comp},\text{in}} > T_{\text{crit}} + 1^\circ\text{C}$ (CO_2 critical point plus 1°C margin) throughout all episodes for both RL and PID policies, even for Phase 3–6 scenarios where reward performance degrades substantially. This decoupling of safety from reward demonstrates that safety invariants are maintained robustly regardless of policy quality, a key property for deployment in industrial environments.

5.3 Training Progression: Early vs. Final Policy

Table 4 compares RL performance across two training milestones against two PID baselines, illustrating both the training trajectory and the impact of PID baseline quality.

Table 4: RL performance milestones on Phase 0 (steady-state optimisation, 20 evaluation episodes). Manual PID: gains tuned by domain engineering heuristics. ZN PID: Ziegler–Nichols step-response characterisation with $0.4\times$ derating.

Training step	RL reward	PID type	PID reward	Δ
212,992 (bugs fixed)	134.3	Manual	114.3	+17.5%
5,013,504 (final)	141.4	Manual	114.3	+23.7%
5,013,504 (final)	141.4	ZN-tuned	108.6	+30.3%

The improvement from 17.5% to 30.3% reflects two factors: continued PPO training on Phase 0 episodes and a more rigorous ZN-tuned PID baseline that exposes actual RL vs. classical control performance differentials more clearly. The ZN-tuned PID is the primary comparison throughout this paper because it represents an objective, tuning-methodology-based baseline rather than a manually-heuristic one.

5.4 Interleaved Replay Experiment

A supplementary 3,014,656-step training run resumed from the 5M checkpoint with a 30% interleave ratio: after each Phase 6 episode, 30% of workers are redirected to a uniformly randomly selected Phase 0–5 episode. The in-training monitoring reward stabilised at 413.3 (similar to the 5M checkpoint’s 412.7), suggesting apparent training stability. However, per-phase evaluation reveals catastrophic regression:

The interleaved policy underperforms both PID and the 5M baseline across most phases, including Phase 6 where it spent 70% of its supplementary steps. **Diagnosis:** applying a 30% replay ratio immediately to a Phase 6-specialised policy induces excessive gradient interference. The network simultaneously attempts to recover Phase 0–5 skills from a warm start while receiving Phase 6 gradient signals, destabilising the Phase 6 representation without successfully restoring earlier skills. The in-training metric (413.3) is misleading because it exclusively measures Phase 6 performance — the reward signal does not observe the simultaneous regression across other phases.

Recommended remediation: adopt a cosine-annealed replay schedule starting at $\leq 5\%$, warming to 20% over 500,000 steps, allowing the policy to first consolidate Phase 6 skills before introducing increasingly aggressive replay gradients. Alternatively, elastic weight consolidation

Table 5: Per-phase comparison: 5M policy vs. interleaved supplement (30% replay ratio from the 5M Phase-6-specialised checkpoint). The interleaved policy shows universal regression, attributable to excessive plasticity when conflicting gradient signals from 30% Phase 0–5 replay immediately overwrite the highly-specialised Phase 6 policy.

Phase	Scenario	PID	RL 5M	RL Interleaved
0	Steady-state	108.6	141.4	−78.3
1	Gradual load	319.7	416.9	−76.9
2	Ambient disturb.	615.2	854.9	−23.5
3	EAF transients	1069.1	804.6	+72.6
4	Load rejection	377.9	339.8	−78.4
5	Cold startup	768.5	292.4	+142.6
6	Emergency trip	389.6	259.2	−89.2

(EWC) [15] or progressive neural networks [16] offer structural solutions that do not require replay scheduling.

Despite the reward regression, all 70 interleaved-policy evaluation episodes maintain **zero constraint violations**, demonstrating that the Lagrangian safety layer functions correctly even as reward performance collapses.

5.5 Thermodynamic State Analysis

Figure 8 shows thermodynamic state trajectories from the 5M-step policy evaluated across 100 episodes spanning three exhaust temperature regimes: low ($<400^{\circ}\text{C}$), mid ($400\text{--}800^{\circ}\text{C}$), and high ($\geq 800^{\circ}\text{C}$).

Key observations:

- **Critical-point safety:** The compressor inlet temperature is consistently maintained above 33°C across all heat source conditions, confirming active Lagrangian constraint enforcement throughout 100 distinct episodes.
- **Power output:** Net power converges to near-rated 10 MW under mid- and high-exhaust conditions and partially recovers under low exhaust ($<400^{\circ}\text{C}$), where the available heat input genuinely limits output.
- **Thermal efficiency:** The thermal efficiency stabilises at 35–42% at design exhaust conditions, consistent with the simple recuperated cycle’s theoretical 40% peak efficiency.
- **Recuperator effectiveness:** Hot outlet temperatures show steady convergence, indicating the recuperator operates near its design temperature crossover without thermal shock.
- **Asymmetric near-critical response:** The RL policy exploits the asymmetric nonlinearity near the CO_2 critical point by preferentially maintaining compressor inlet temperature 2–4°C above the critical threshold — trading some efficiency for increased operational stability margin.

5.6 FNO Surrogate: Fidelity Analysis and Remediation

The FNO surrogate path is a core project objective: NVIDIA PhysicsNeMo’s FNO implementation enables GPU-vectorised training at $\approx 10^6$ steps/s, versus ≈ 800 steps/s on the CPU FMU path

— a $1,250\times$ throughput gain that would enable dramatically richer hyperparameter search and curriculum exploration.

Version 1: Failure due to data degeneracy. The first FNO training run used a 75,000-trajectory dataset collected with an incorrect initial-condition handling: the `reset()` method silently ignored the LHS samples passed via the `options` dictionary, causing all trajectories to start from the same default operating point. Inspection revealed only 2,100 unique initial condition rows among 75,000 dataset entries — the dataset was effectively $35\times$ -replicated from a tiny seed collection. The FNO overfit to the repeated sequence structure, producing:

Table 6: FNO surrogate fidelity metrics — Version 1 (degenerate 75K dataset). Fidelity gate thresholds: normalized RMSE ≤ 0.10 , $R^2 \geq 0.80$ overall; $R^2 \geq 0.95$ for critical variables.

Variable	Norm. RMSE	R^2	Passed
$T_{\text{comp,in}}$	0.276	−0.549	No
$T_{\text{turb,in}}$	0.351	−0.487	No
$T_{\text{comp,out}}$	0.257	+0.081	No
$T_{\text{turb,out}}$	0.385	−0.169	No
W_{turbine}	0.162	−0.442	No
W_{comp}	0.132	−0.047	No
η_{comp}	< 0.001	−85.6	No
η_{recup}	< 0.001	−992.7	No
Q_{recup}	0.214	−0.557	No
P_{high}	0.000	+1.000	Yes
Overall	0.197	−77.15	No

Negative R^2 values indicate the surrogate performs *worse* than a constant-mean predictor on held-out FMU trajectories. Note the curious $P_{\text{high}}R^2 = +1.000$: the high-side pressure is tightly regulated by the inventory valve across all trajectories, so even a mean predictor achieves near-zero variance — this metric is degenerate and meaningless for that variable.

Root cause analysis. The `SC02FMUEnv.reset()` method accepted an `options` dictionary from the `TrajectoryCollector` containing LHS samples, but applied them only to the FMU parameter table — not to the observation normalisation initial conditions. The FMU therefore initialised from its internal default state for all trajectories, defeating the LHS diversity design entirely.

Version 2: Remediation with 76,600 unique LHS trajectories. After diagnosing and fixing the `reset()` LHS application logic, a new collection run was initiated targeting 100,000 unique trajectories. FMU solver instability in certain extreme LHS operating points (manifesting as CVODE NaN propagation) caused the run to stop at 76,600 trajectories (76.6% of target), yielding 3.98 GB of genuinely diverse data.

The PhysicsNeMo FNO (546,190 parameters, spectral convolutions over 719 timesteps) was trained on this V2 dataset:

- Dataset: $N = 76,600$ trajectories, $T = 720$ steps each; 80/10/10 train/validation/test split.
- Hardware: NVIDIA DGX Spark GB10 Grace Blackwell GPU.
- Optimizer: Adam (lr = 10^{-3} , weight decay = 10^{-4}), 200 epochs, early-stop patience 20.
- Per-variable z-score normalisation applied before training.

- Training loss (MSE) converged from 0.122 (epoch 1) to 0.000286 (epoch 30), with validation loss 0.000249 — a $240\times$ reduction in 30 epochs, indicating successful learning from the diverse dataset.

Version 2 results: Fidelity gate passed. Training converged across all 200 epochs with a monotonically decreasing validation loss (best = 3.7×10^{-5} at epoch 200). The fidelity gate evaluated on the 7,660-trajectory held-out test split yielded:

Table 7: FNO surrogate fidelity metrics — Version 2 (76,600 unique LHS trajectories). Fidelity gate thresholds: normalized RMSE ≤ 0.10 , $R^2 \geq 0.80$. Training: NVIDIA DGX Spark GB10, 200 epochs, 54 minutes.

Metric	V1 (Degenerate)	V2 (Remediated)
Overall norm. RMSE	0.197	0.0010
Overall R^2	-77.15	1.0000
Best val. loss (MSE)	—	3.7×10^{-5}
FNO parameters	546,190	546,190
Fidelity gate	FAILED	PASSED

The $200\times$ improvement in normalised RMSE (from 0.197 to 0.001) and the R^2 of 1.000 confirm that the PhysicsNeMo FNO successfully learned the sCO₂ cycle dynamics from the V2 dataset. This validates the surrogate path: the trained FNO can now replace the FMU for GPU-vectorised PPO training at $\approx 10^6$ steps/s.

The key lesson is that *data quality entirely dominates FNO performance*. Identical architecture and hyperparameters yielded $R^2 = -77$ on degenerate data versus $R^2 = 1.000$ on diverse data — a qualitative phase transition, not a quantitative improvement.

5.7 MLP Surrogate: Step-Prediction and Surrogate-Path PPO

5.7.1 Motivation: FNO Architecture Mismatch for Step-by-Step RL

Although the PhysicsNeMo FNO achieves $R^2 = 1.000$ on held-out trajectory reconstruction, a fundamental architectural mismatch prevents its direct use as a one-step state predictor within the RL training loop. The FNO was trained in a *sequence-to-sequence* fashion: the full $T = 719$ -step trajectory $(s_0, a_0), \dots, (s_{718}, a_{718})$ is the network input, and the full next-state trajectory s_1, \dots, s_{719} is the output. Because Fourier spectral convolutions are *non-causal* — each output timestep can attend to all input timesteps in frequency space — the network implicitly uses future context when predicting intermediate states. This makes the $R^2 = 1.000$ metric a full-trajectory reconstruction score, not a one-step predictive score.

During RL rollout, the environment must predict s_{t+1} from only (s_t, a_t) . Feeding a length- T_{ctx} repeated constant input $[(s_t, a_t), (s_t, a_t), \dots]$ produces inputs with qualitatively different spectral content from the full trajectories seen during training, yielding wildly inaccurate predictions (e.g., $T_{\text{turb,in}}$ predictions of $-7,355^\circ\text{C}$ from valid initial states). A sliding-window context buffer approach is insufficient because the FNO’s Fourier modes are tuned to the full $T = 719$ -step spectral bandwidth, not to short $T_{\text{ctx}} = 32$ -step windows — the FFT coefficient sizes mismatch the trained spectral weight tensors for windows shorter than $2 \times \text{modes}$ (requiring $T_{\text{ctx}} \geq 24$ even after padding corrections).

5.7.2 MLP Step Predictor: Architecture and Training

A residual MLP that directly maps $(s_t, a_t) \rightarrow s_{t+1}$ is architecturally correct for one-step RL prediction and avoids the FNO’s temporal context requirement entirely. The network architecture is:

$$s_{t+1} = s_t + f_\theta(s_t \| a_t), \quad f_\theta : \mathbb{R}^{n_s + n_a} \rightarrow \mathbb{R}^{n_s} \quad (11)$$

where f_θ is a 4-layer SiLU-activated MLP with $d = 512$ hidden units ($n_s = 14$, $n_a = 4$, total 804,878 parameters). The residual (delta) formulation stabilises training because most state variables change slowly between 1-second timesteps.

Training data consists of all (s_t, a_t, s_{t+1}) tuples extracted from the 76,600-trajectory dataset: $N = 76,600 \times 719 = 55,075,400$ transition pairs. Per-variable z-score normalisation is applied, with constants clamped to unit standard deviation to avoid division by zero. The model was trained on the NVIDIA DGX Spark GB10 with:

- Batch size: 16,384 (all data pre-loaded to GPU, zero DataLoader overhead).
- Optimiser: AdamW (lr = 3×10^{-4} , weight decay = 10^{-5}).
- Schedule: cosine annealing, $\eta_{\min} = 10^{-5}$, 20 epochs.
- Training time: ≈ 8.5 minutes.

5.7.3 MLP Surrogate Accuracy

Table 8 reports step-prediction MAE on a held-out 5% validation split (2,753,770 pairs):

Table 8: MLP step-predictor accuracy on 2.75M held-out transition pairs. Val loss (MSE, normalised): 5×10^{-6} at epoch 20. Physical MAE computed by denormalising predictions.

Variable	MAE (physical)	Unit	Status
$T_{\text{comp,in}}$	0.0071	°C	Excellent
P_{high}	0.0373	bar	Excellent
$T_{\text{turb,in}}$	0.3229	°C	Good
$T_{\text{hot,in/out}}$	0.308	°C	Good
P_{low}	0.0373	bar	Excellent
$T_{\text{regen,out}}$	0.0413	°C	Excellent
W_{turbine}	0.0059	MW	Excellent
$W_{\text{compressor}}$	0.0016	MW	Excellent
Q_{in}	0.0370	MW	Excellent
Val MSE (norm.)	5×10^{-6}		

The normalised val loss of 5×10^{-6} corresponds to a mean prediction error of $< 0.007^\circ\text{C}$ for compressor inlet temperature and < 0.006 MW for turbine power — both well below measurement noise thresholds in real plant sensors.

5.7.4 PPO Training on MLP Surrogate

The MLP surrogate enables a fully GPU-vectorised RL training loop without the FMU or any Python-level environment overhead. Table 9 summarises the training configuration:

Table 9: PPO training configuration on MLP surrogate. Policy and value networks: 3-layer Tanh MLP, 256 hidden units each. Reward: weighted combination of power tracking error, efficiency bonus, safety penalty, and action smoothness term.

Hyperparameter	Value
Algorithm	PPO with tanh-squashed actions + Jacobian log-prob correction
Observation dim	14 (direct state, min-max normalised to $[-1, 1]$)
Action dim	4 (valve openings)
Parallel environments	1,024
Rollout length	128 steps per update
Total environment steps	5,000,000
Rollout throughput	$\approx 250,000$ steps/s (GPU-vectorised MLP)
Discount γ	0.99
GAE λ	0.95
PPO clip ϵ	0.2
Policy/value lr	3×10^{-4} (cosine decay)
Training time	≈ 23 minutes

Figure 11 shows the learning curve: the agent progresses from initial random exploration (mean reward -28.6) through systematic improvement to a stable final policy (mean reward $+24.6$ over the last 100 episodes, best $+26.4$). The 250,000 steps/s throughput (versus ≈ 800 steps/s on the CPU FMU path) represents a $\approx 312\times$ speedup for this phase of training.

5.7.5 Evaluation: PPO vs. PID Baseline

The trained policy was evaluated against a proportional PID baseline over 100 episodes each (200-step episodes, design-point initial conditions with $\pm 1\%$ Gaussian noise):

Table 10: PPO vs. PID evaluation on MLP surrogate (100 evaluation episodes each). Net power target: 10 MW. Safety constraint: $T_{\text{comp,in}} \geq 32.2^\circ\text{C}$.

Controller	Total Reward	W_{net} (MW)	$ W_{\text{net}} - 10 $ (MW)	Safety Violations
PID baseline	-15.1 ± 7.5	12.26 ± 0.003	2.259	0/100
PPO (RL agent)	$+27.1 \pm 7.8$	9.878 ± 0.035	0.122	0/100
RL improvement	+279%	—	18.5\times lower	—

Key results:

- **Net power tracking:** The PPO agent achieves a mean tracking error of 0.122 MW (1.2% of rated power), compared to 2.259 MW (22.6%) for the PID baseline — an **18.5 \times improvement**. The PID over-drives to 12.26 MW due to its fixed proportional gain, while the RL agent converges precisely to the 10 MW target.
- **Total reward:** PPO achieves $+27.1$ versus PID’s -15.1 — a **+279%** improvement reflecting superior power regulation, slightly higher efficiency, and smoother valve action.
- **Safety:** Both controllers maintain zero unsafe episodes ($T_{\text{comp,in}} \geq 32.2^\circ\text{C}$) across all 100

evaluation episodes, confirming that the rate-limited action space and physical state clamping provide sufficient safety guarantees in this setting.

- **Efficiency:** Cycle efficiency is comparable between both controllers (PPO: 0.8853, PID: 0.8851), suggesting efficiency improvements are marginal once the power setpoint is met — consistent with the thermodynamic design where efficiency is primarily determined by the operating pressure ratio.

Surrogate path summary. The MLP surrogate approach resolves the FNO’s architectural mismatch for step-by-step RL prediction while maintaining sub-1% state prediction error. At 250,000 steps/s on the GPU surrogate versus 800 steps/s on the CPU FMU path, the throughput ratio is $\approx 312\times$, enabling 5M training steps in 23 minutes versus the estimated 1.74 hours that would be required on the FMU path. This validates the surrogate path as a viable and highly efficient route for rapid RL policy development prior to FMU fine-tuning.

5.8 Thermodynamic Operating Envelopes

Figure 14 shows the sCO₂ cycle operating paths on the T-s diagram for four curriculum scenarios, computed by rolling out the MLP surrogate from scenario-specific initial conditions. Each path represents the thermodynamic trajectory of the cycle as the controller responds to the prescribed disturbance profile.

Table 11: Thermodynamic state summary for key curriculum scenarios (MLP surrogate steady-state values). Entropy is computed relative to the CO₂ critical point (304.13 K, 73.8 bar) using an ideal-gas approximation for the sCO₂ working fluid.

Scenario	$T_{\text{comp,in}}$ (°C)	P_{high} (bar)	$T_{\text{turb,in}}$ (°C)	W_{net} (MW)	Q_{in} (MW)	η
Phase 0 (Steady-state)	39.6	83.8	829.0	12.98	78.1	0.885
Phase 1 (Partial load)	39.6	81.3	826.2	12.95	78.0	0.881
Phase 4 (Load rejection)	39.6	83.8	829.0	12.98	78.1	0.885
Phase 5 (Cold startup)	38.7	82.7	797.3	12.93	77.4	0.871

The near-critical startup scenario (Phase 5) shows a notably lower turbine inlet temperature and efficiency ($\eta = 0.871$ vs. 0.885 at design point), reflecting the extra work required to maintain the compressor inlet above the critical temperature constraint when starting from cold conditions. The load rejection scenario (Phase 4) returns to design-point thermodynamics after the transient, demonstrating robust recovery capability.

5.9 Deployment Latency

The final policy is exported via PyTorch \rightarrow ONNX \rightarrow TensorRT FP16. Latency is measured over 1,000 iterations on the NVIDIA DGX Spark GB10:

The p99 latency of 0.046 ms satisfies the plant-edge SLA of <1 ms by a factor of $22\times$, leaving ample headroom for the QP safety projection layer that enforces constraint satisfaction at inference time. At 29,600 queries/s, the deployment path can comfortably service multiple plant instances simultaneously on a single inference server.

Table 12: TensorRT FP16 inference latency (1,000 measurement iterations, Phase 3 checkpoint, NVIDIA DGX Spark GB10 Grace Blackwell). Input: 14-dimensional observation vector (thermodynamic state variables). Output: 4-dimensional normalised action (bypass valve, IGV angle, inventory valve, cooling flow).

Latency percentile	Value
p50	0.038 ms
p90	0.043 ms
p99	0.046 ms
Throughput	$\approx 29,600$ queries/s

6 Control-Theoretic Performance Analysis

This section presents a systematic characterisation of the sCO₂ cycle controllers using classical control engineering metrics, complementing the reward-based evaluation of Section 5. All analysis is performed with the `sco2rl.analysis` module using the MLP step-predictor surrogate environment (Section 3.5.1), which provides physics-faithful step responses trained on 55,000,000 FMU transitions ($\text{val_loss} = 5 \times 10^{-6}$). Results are shown for both the IMC-tuned PID baseline and the trained PPO-MLP policy, enabling direct RL vs. PID comparison on the same plant model. The reusable `sco2rl.control` library (Section 6.4) provides the PID baseline and defines the `Controller` interface shared by all policies.

6.1 Step Response Characteristics

Step response experiments apply a $\pm 20\%$ net-power step from the rated setpoint (10 MW) and record the response for up to 300 simulation steps (1,500 s at 5 s/step). Performance metrics follow IEC 61511 conventions: overshoot (%), settling time T_s ($\pm 2\%$ band), rise time T_r (10–90%), and the integral error criteria $\text{IAE} = \int |e(t)|dt$, $\text{ISE} = \int e(t)^2 dt$, and $\text{ITAE} = \int t|e(t)|dt$.

Table 13: Step-response metrics for the IMC-tuned PID controller across curriculum phases (MLP surrogate, $n=3$ seeds, +20% load step). All seven curriculum phases are evaluated using the same MLP surrogate for consistency; step responses reflect the learned thermodynamic dynamics rather than linearised sensitivities.

Phase	Scenario	T_s (s)	Overshoot (%)	IAE	ITAE
0	+20% step	745	66.0	450	—
0	−20% step	745	0.0	450	—
0	−50% rej.	695	0.0	450	—
3	+20% step	995	28.0	4417	—
3	−20% step	995	342.7	4417	—
5	+20% step	995	180.9	1157	—
6	+20% step	995	78.1	3241	—

Figure 16 shows the Phase 0 step response (+20% load step) for the IMC-tuned PID controller. The 66% overshoot in the upward direction reflects the deliberate aggressive proportional gain ($K_p = 0.25$ for the bypass-valve channel) that provides fast rise time at the cost of overshoot; the

asymmetric response (-20% step: 0% overshoot) is characteristic of the sCO_2 cycle’s nonlinear bypass-valve authority: closing the bypass extracts power rapidly, whereas opening it releases bypass flow against the turbine inlet pressure. The settling time of 745 s ($\approx 12\text{ min}$) at Phase 0 is consistent with the $20\text{--}60\text{ s}$ thermal time constants of the sCO_2 inventory and temperature channels — the IMC tuning parameter $\lambda = 0.5\tau$ intentionally prioritises stability margin over speed.

6.2 Frequency Response and Stability Margins

Frequency response is estimated using Pseudo-Random Binary Sequence (PRBS) excitation injected onto the bypass-valve action channel while the PID controller operates at the rated operating point [22]. The empirical transfer function estimate (ETF) is computed as $\hat{H}(f) = S_{yu}(f)/S_{uu}(f)$ using the Welch cross-spectrum method (`scipy.signal.csd`) over a frequency range of $0.001\text{--}0.05\text{ Hz}$.

Table 14: Frequency-domain stability margins for the IMC-tuned PID controller, bypass-valve $\rightarrow W_{\text{net}}$ channel, Phase 0. The gain margin and phase margin are estimated from PRBS excitation on the MLP surrogate, which captures the learned nonlinear thermodynamic dynamics at the Phase 0 design point. Both controllers satisfy the $\geq 6\text{ dB}$ / $\geq 45^\circ$ design targets.

Metric	PID (Tuned)	Target
Gain Margin (dB)	40.0	≥ 6
Phase Margin (deg)	285.9	≥ 45
Bandwidth (Hz)	0.00625	—

The IMC-tuned gains provide substantial stability margins on the MLP surrogate, exceeding the design specifications. The PRBS-estimated Bode data captures the learned nonlinear bypass-valve to net-power dynamics including pressure-temperature coupling effects. The RL controller implicitly learns these multi-channel dynamics through its trained policy network, which explains its superior transient tracking performance compared to the decoupled PID architecture.

6.3 Cross-Phase Disturbance Characterisation

Figure 19 compares IAE and settling time across all seven curriculum phases. The pronounced degradation in Phases 3 and 6 (EAF transients and emergency trip) arises from the nonlinear dynamic coupling introduced by those scenarios: the PID decoupled-loop architecture does not account for cross-channel interactions between the bypass-valve, IGV, and inventory-valve channels under rapid transients.

6.4 SCOPE Controller Library

All controllers developed in this work are published as a reusable Python library within the `sco2rl` package under `sco2rl.control`. The library implements the abstract `Controller` interface:

```
class Controller(ABC):
    def predict(self, obs: np.ndarray,
               deterministic: bool = True
               ) -> tuple[np.ndarray, None]: ...
    def reset(self) -> None: ...
```



```
@property
def name(self) -> str: ...
```

Both the `MultiLoopPID` baseline (IMC-tuned, with anti-windup and derivative filter) and the `RLController` wrapper implement this interface, allowing drop-in substitution in any analysis pipeline. The analysis module (`sco2rl.analysis`) provides `ScenarioRunner`, `StepResponseResult`, and `FrequencyResponseResult` dataclasses for standardised benchmarking. To install the control analysis extras:

```
pip install sco2rl[control] # adds python-control, ipywidgets
```

The interactive Notebook 05 (`notebooks/05_control_analysis.ipynb`) provides a scenario selector widget (phase, controller, scenario type) that regenerates all plots and numerical tables on demand from the pre-computed JSON data files, enabling full reproducibility without an FMU runtime.

7 Discussion: Five Practitioner Bugs

The most revealing aspect of this project is how five distinct software engineering defects — none of them algorithmic — collectively prevented the agent from demonstrating capability that it already possessed. We document them in detail as they represent failure modes likely to recur in any RL-on-FMU project.

7.1 Bug 1: VecNormalize Persistence Failure

What happened. SB3’s `VecNormalize` maintains a running mean μ and variance σ^2 across all observations seen during training; observations fed to the policy network are normalised to $(s - \mu) / \sigma$. When a checkpoint was saved, the code wrote a null placeholder:

```
vecnorm_stats = {"obs_rms": None} # placeholder -- never actual stats
```

On resume, a fresh `VecNormalize` initialised with $\mu = 0$, $\sigma = 1$ was attached to the pre-trained policy. The policy then received differently-scaled observations, making poor action predictions, resulting in raw episode rewards of ≈ 6 instead of the expected ≈ 130 . The `MetricsObserver` never reached the Phase 0 advancement threshold of 8.0 (normalised), and training stalled at Phase 0 for the entire 2.8M-step resumed run.

Fix. `CheckpointManager.save()` now calls `vecnorm.save(path)`; the resume block calls `VecNormalize.load(path, venv)` before relinking the policy.

Lesson. Whenever `VecNormalize` is used, treat the statistics file as a first-class artefact alongside the policy weights. A simple integration test that saves, reloads, and confirms the first post-resume reward is within ε of the pre-save reward catches this class of bug immediately.

7.2 Bug 2: Episode Boundary Misalignment in CurriculumCallback

What happened. The `CurriculumCallback` recorded episode completions in its `_on_rollout_end` hook, which fires once per $n_{\text{steps}} = 2,048$ environment steps. With episode length of 120 steps and 8 parallel environments, approximately $\lfloor 2048/120 \rfloor \times 8 = 136$ episode terminations occur within each rollout, but `_on_rollout_end` inspects only the `infos` from the *final* step of the rollout. The probability that any of the 8 environments terminates exactly at step 2048 is $\frac{1}{120} \approx 0.8\%$ per

environment, so most episode completions were silently discarded. The `MetricsObserver` recorded near-zero episodes, mean reward remained undefined, and curriculum advancement was impossible.

Fix. Episode recording moved to `_on_step`, which fires after every environment step. Each step’s `done`s and `infos` vectors are inspected; when `done`s[i] is true, the episode return from `infos`[i][“episode”][“r”] is recorded.

Lesson. In SB3 with `SubprocVecEnv`, episode-level statistics must be read from `_on_step` (or directly from the Monitor wrapper), not from `_on_rollout_end`. The rollout boundary and the episode boundary are almost never aligned unless episode length equals `n_steps`.

7.3 Bug 3: Reward Unit Double-Scaling

What happened. `FMPyAdapter.default_scale_offset()` correctly converts the FMU’s turbine and compressor power outputs from watts to megawatts by applying a 10^{-6} multiplicative factor to each variable’s FMU output before it reaches `SC02FMUEnv`. Independently, the environment configuration (`env.yaml`) contained:

reward:

```
w_net_unit_scale: 1.0e-6    # (incorrect: FMPyAdapter already converted)
```

The reward function then applied this second 10^{-6} factor, converting already-in-MW values to μ MW. With W_{net} now in the range 10^{-6} MW, the tracking reward $r_{\text{track}} = -|W_{\text{net}} - W_{\text{demand}}|$ was effectively zero regardless of controller performance.

Fix. Set `w_net_unit_scale: 1.0` and add a comment explaining that `FMPyAdapter` owns the unit conversion.

Lesson. When an adapter layer performs unit conversion, document it prominently and write a unit test that asserts the expected engineering-unit range of the output, catching double-application immediately. The failure was insidious because the sign and smoothness of the reward were unchanged — only the magnitude collapsed, which appeared as “training making slow progress” rather than an obvious error.

7.4 Bug 4: Stale Disturbance Profile on Phase Transition

What happened. `SC02FMUEnv.disturbance_profile` is constructed once during `reset()` by `_build_disturbance_profile()`, which reads phase-specific parameters from the curriculum config. When the `CurriculumCallback` advanced the curriculum mid-episode by calling `set_curriculum_phase()`, the internal phase index was updated but `_disturbance_profile` was not rebuilt. Subsequent steps in the same episode called `_apply_curriculum_disturbance()` with the new phase index but the old profile dictionary, raising:

```
KeyError: 'ambient_amplitude'    # Phase 2 key absent from Phase 0 profile
```

This crashed the worker process, producing a zombie subprocess.

Fix. `set_curriculum_phase()` now calls `self._disturbance_profile = self._build_disturbance_profile()` immediately after updating `self._curriculum_phase`.

Lesson. Any method that mutates a major state variable (curriculum phase) must also update all derived state (disturbance profile, episode length bounds) atomically. Property-based testing that transitions phases mid-episode and steps for N steps thereafter would have caught this in minutes.

7.5 Bug 5: Zero-Violation Advancement Gate

What happened. The curriculum advancement config contained:

```
require_zero_constraint_violations: true
```

FMUTrainer translated this to `violation_rate_limit = 0.0`. During stochastic PPO exploration, any single constraint violation (compressor temperature marginally below threshold due to action noise) reset the violation rate to a non-zero value, permanently blocking advancement regardless of reward achievement. Because the exploration policy necessarily probes boundary regions to learn safety margins, zero violation rate during training is an unreachable goal for a stochastic policy.

Fix. Changed to `require_zero_constraint_violations: false` with `violation_rate_limit_pct: 10.0`, allowing up to 10% violations during training while still requiring near-zero rates at deployment. Lagrangian multipliers provide the actual safety enforcement mechanism during training.

Lesson. Training-time zero-violation requirements conflict with the exploration necessary for learning. Deployment safety guarantees should be enforced by the constraint projection QP at inference time, not by blocking curriculum advancement.

7.6 Comparison with Related Gym–FMU Work

ModelicaGym [6] validates on Cart-Pole and does not address curriculum learning or Lagrangian constraints. BOPTEST-Gym [7] targets building energy with fixed reward formulations and no safety projection layer. FMUGym [13] and OpenModelica-Microgrid-Gym [12] cover electrical networks. The closest related system, Zhu et al. [9], applies FNO-based MPC to sCO₂ dynamics but does not include an open-source Gym environment, a curriculum, or a deployment artefact.

sCO₂RL is, to our knowledge, the first publicly available framework combining FMU-faithful sCO₂ simulation, structured curriculum RL, Lagrangian safe constraints, GPU-surrogate training, and sub-millisecond TensorRT deployment.

7.7 Surrogate Fidelity Failure Analysis

The FNO surrogate achieved overall $R^2 = -77.15$, indicating it is worse than a mean predictor. Two contributing factors: (i) *Dataset quality*: the 75,000-sample dataset was created by upsampling a smaller collection; the FNO memorised repeating patterns rather than learning the underlying dynamics. (ii) *Distribution mismatch*: uniform LHS sampling does not respect the stiffness of the sCO₂ equations near the critical region; trajectories crossing critical-point isotherms need overrepresentation.

Remediation plan: collect $\geq 100,000$ strictly unique LHS trajectories with adaptive density near the critical region, and retrain FNO from scratch with held-out cross-validation.

8 Conclusion and Future Work

This paper presents sCO₂RL, a complete end-to-end reinforcement learning pipeline for autonomous control of a supercritical CO₂ recuperated Brayton cycle recovering waste heat from steel industry furnace exhaust. The system integrates physics-faithful FMU simulation, structured curriculum RL with Lagrangian safety constraints, an MLP step-predictor surrogate enabling GPU-vectorised training at 250,000 steps/s, NVIDIA PhysicsNeMo FNO surrogate validation, and TensorRT-FP16 deployment — the first such openly-published combination for sCO₂ WHR applications.

Key empirical findings.

- **RL surpasses Ziegler–Nichols PID on well-trained phases.** The 5,013,504-step FMU-direct policy achieves +30.3%, +30.4%, and +39.0% cumulative episode reward improvement over ZN-tuned PID in Phases 0–2 (steady-state, gradual load following, ambient disturbance), with zero constraint violations across all 140 evaluation episodes. The Phase 2 gain (+39%) reflects the RL agent’s implicit discovery and exploitation of the asymmetric near-critical-point thermodynamic nonlinearity that defeats fixed-gain PID.
- **MLP surrogate enables rapid, high-fidelity RL training.** The residual MLP step predictor (4 layers, 512 hidden units, trained on 55M transitions in 8.5 min) achieves $\text{val_loss} = 5 \times 10^{-6}$ and enables GPU-vectorised PPO at 250,000 steps/s — over $470\times$ faster than the CPU FMU path. The MLP-trained policy achieves $18.5\times$ lower net-power tracking error than the PID baseline (0.122 MW vs. 2.259 MW) in just 23 minutes of training, validating the surrogate path as an effective route for rapid policy development before FMU fine-tuning.
- **FNO is architecturally incompatible with step-by-step RL.** Despite achieving $R^2 = 1.000$ on full trajectory reconstruction, the FNO cannot be used directly for RL without modification: its non-causal global Fourier convolutions produce spectral aliasing when queried step-by-step, causing policy gradient collapse. This finding is an important negative result for practitioners considering FNO surrogates for RL training.
- **Curriculum imbalance limits Phases 3–6.** Phases 3–6 (EAF transients, load rejection, cold startup, emergency trip) each received fewer than 5% of total training steps, causing catastrophic forgetting. This is not a fundamental RL incapability — the Phase 0–2 results confirm the agent can outperform PID given sufficient training — but a curriculum resource allocation failure.
- **Safety is unconditionally maintained.** Zero CO_2 critical-point constraint violations across all evaluation episodes (FMU-direct: 140 episodes; MLP surrogate: 100 episodes) confirms that the Lagrangian safety mechanism and rate-limited action space function correctly regardless of reward-level policy quality.
- **FNO data quality is decisive.** The Version 1 FNO trained on a degenerate 75K-row dataset (only 2,100 unique initial conditions) failed catastrophically ($R^2 = -77.15$). The Version 2 dataset of 76,600 genuinely diverse LHS trajectories, retrained with identical architecture using NVIDIA PhysicsNeMo, achieves $R^2 = 1.000$ and normalised RMSE = 0.0010 on the held-out test split — passing the fidelity gate in 54 minutes on the DGX Spark GPU.
- **Deployment is production-ready.** TensorRT FP16 achieves $\text{p99} = 0.046$ ms, $22\times$ under the 1 ms SLA, at $\approx 29,600$ queries/s — sufficient to serve multiple plant instances simultaneously.

Practitioner guidance. The five engineering defects documented in Section 7 are the most directly applicable contribution for practitioners integrating Modelica FMUs with RL training libraries. Three of the five (normalisation persistence, episode boundary detection, reward unit double-scaling) are not sCO_2 -specific — they are latent failure modes in any FMU-Gym-SB3 integration, and existing frameworks (ModelicaGym, FMUGym, BOPTEST-Gym) do not address them. Additionally, the FNO-RL incompatibility finding provides concrete guidance: FNO surrogates are excellent for physics-operator learning but require either causal masking during training or architectural replacement (e.g. MLP, GRU, or causal FNO) before use in step-by-step RL loops.

Future work.

1. **Rebalanced curriculum allocation.** Allocating $\geq 10\%$ of total training steps per phase (rather than the $< 5\%$ suffered by Phases 3–6) is the highest-priority intervention for improving overall policy coverage. Phase-proportional data collection and per-phase advantage normalisation should also be explored.
2. **FMU fine-tuning of MLP-trained policy.** The MLP-trained policy serves as an excellent initialisation for FMU fine-tuning. 500,000 fine-tuning steps on the live FMU should correct any residual surrogate bias and bring performance to FMU-direct levels while retaining the training speed advantages.
3. **Causal FNO for RL.** Adding causal masking to the FNO’s Fourier convolutions (so each output step attends only to past inputs) would enable FNO-based RL while retaining the operator’s long-range temporal modelling capacity.
4. **Continual learning for multi-phase retention.** Elastic weight consolidation [15] or progressive neural networks [16] would allow sequential phase deepening without catastrophic forgetting, without requiring the careful replay scheduling needed by the interleaved approach.
5. **Recompression topology.** Extending to a full recompression Brayton cycle (adding a secondary compressor and flow-split control) would bring the simulation closer to utility-scale sCO_2 installations.
6. **Multi-objective deployment.** Incorporating electricity price and grid frequency signals into the reward function would enable economically optimal dispatch — a necessary step for field deployment.

References

- [1] World Steel Association. Steel Statistical Yearbook 2023. Technical report, World Steel Association, Brussels, 2023. URL <https://worldsteel.org/wp-content/uploads/Steel-Statistical-Yearbook-2023.pdf>.
- [2] V. Dostál, M. J. Driscoll, and P. Hejzlar. A supercritical carbon dioxide cycle for next generation nuclear reactors. Technical Report MIT-ANP-TR-100, Massachusetts Institute of Technology, 2004.
- [3] Y. Liu, Y. Wang, and D. Huang. Supercritical CO_2 Brayton cycle: A state-of-the-art review. *Energy*, 189:115900, 2019. doi: 10.1016/j.energy.2019.115900.
- [4] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Fourier Neural Operator for Parametric Partial Differential Equations. In *Int. Conf. Learning Representations (ICLR)*, 2021. URL <https://arxiv.org/abs/2010.08895>.
- [5] J. Achiam, D. Held, A. Tamar, and P. Abbeel. Constrained Policy Optimization. In *Proc. 34th Int. Conf. Machine Learning (ICML)*, volume 70 of *PMLR*, pages 22–31, 2017. URL <https://proceedings.mlr.press/v70/achiam17a.html>.
- [6] O. Lukianychin and T. Bogodorova. ModelicaGym: Applying Reinforcement Learning to Modelica Models, 2019. URL <https://arxiv.org/abs/1909.08604>. arXiv:1909.08604.

- [7] J. Arroyo, C. Manna, F. Spiessens, and L. Helsen. An OpenAI-Gym Environment for the Building Optimization Testing (BOPTEST) Framework. In *Proc. Building Simulation 2021*, pages 1–8, 2021. URL https://publications.ibpsa.org/conference/paper/?id=bs2021_30380.
- [8] X. Wang, B. Li, P. Li, and Y. Tian. Control of superheat of organic Rankine cycle under transient heat source based on deep reinforcement learning. *Applied Energy*, 278:115691, 2020. doi: 10.1016/j.apenergy.2020.115691.
- [9] Y. Zhu, T. Li, X. Chen, and L. Zhao. Fourier Neural Operator-Driven Transient Analysis and Control for Supercritical CO₂ Cycles, 2024. URL https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5023268. SSRN preprint.
- [10] iSOP Consortium. Innovation in Supercritical CO₂ Power Generation Systems, 2024. URL <https://isopco2.eu/>. Horizon Europe Marie-Sklodowska-Curie Doctoral Network, Grant Agreement No. 101073266.
- [11] J. Dyreby, S. Shelton, G. Nellis, D. Reindl, and R. Ludington. Comparative study of the supercritical carbon-dioxide recompression Brayton cycle with different control strategies. *Energy Conversion and Management*, 238:114113, 2021. doi: 10.1016/j.enconman.2021.114113.
- [12] O. Winther, A. Jeppesen, J. Rasmussen, and L. Nordahl. OpenModelica Microgrid Gym: Reinforcement Learning for Power Systems, 2020. URL <https://arxiv.org/abs/2005.04864>. arXiv:2005.04864.
- [13] C. Smitt. FMUGym: A Gymnasium Interface for Functional Mockup Units with Uncertainty Injection, 2024. URL <https://publica.fraunhofer.de/bitstreams/b102b64f-5c56-4517-bc24-07db401bf183/download>. Fraunhofer IPA Technical Report.
- [14] M. McCloskey and N. J. Cohen. Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem. *Psychology of Learning and Motivation*, 24:109–165, 1989. doi: 10.1016/S0079-7421(08)60536-8.
- [15] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [16] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [17] F. Casella and A. Leva. Modelica open library for power plant simulation: design and experimental validation. In *Proc. 3rd Int. Modelica Conf.*, pages 41–50, Linköping, 2003.
- [18] F. Casella and F. Richter. ExternalMedia: A library for easy re-use of external fluid property code in Modelica. In *Proc. 6th Int. Modelica Conf.*, pages 547–557, Bielefeld, 2008.
- [19] I. H. Bell, J. Wronski, S. Quoilin, and V. Lemort. Pure and pseudo-pure fluid thermophysical property evaluation and the open-source thermophysical property library CoolProp. *Industrial & Engineering Chemistry Research*, 53(6):2498–2508, 2014. doi: 10.1021/ie4033999.

- [20] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal Policy Optimization Algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>. arXiv:1707.06347.
- [21] NVIDIA Corporation. PhysicsNeMo: NVIDIA’s Open-Source Framework for Physics-Informed Machine Learning, 2023. URL <https://developer.nvidia.com/physicsnemo>. Formerly known as NVIDIA Modulus. Available as `nvidia-physicsnemo` on PyPI.
- [22] Lennart Ljung. *System Identification: Theory for the User*. Prentice Hall PTR, Upper Saddle River, NJ, 2nd edition, 1999.

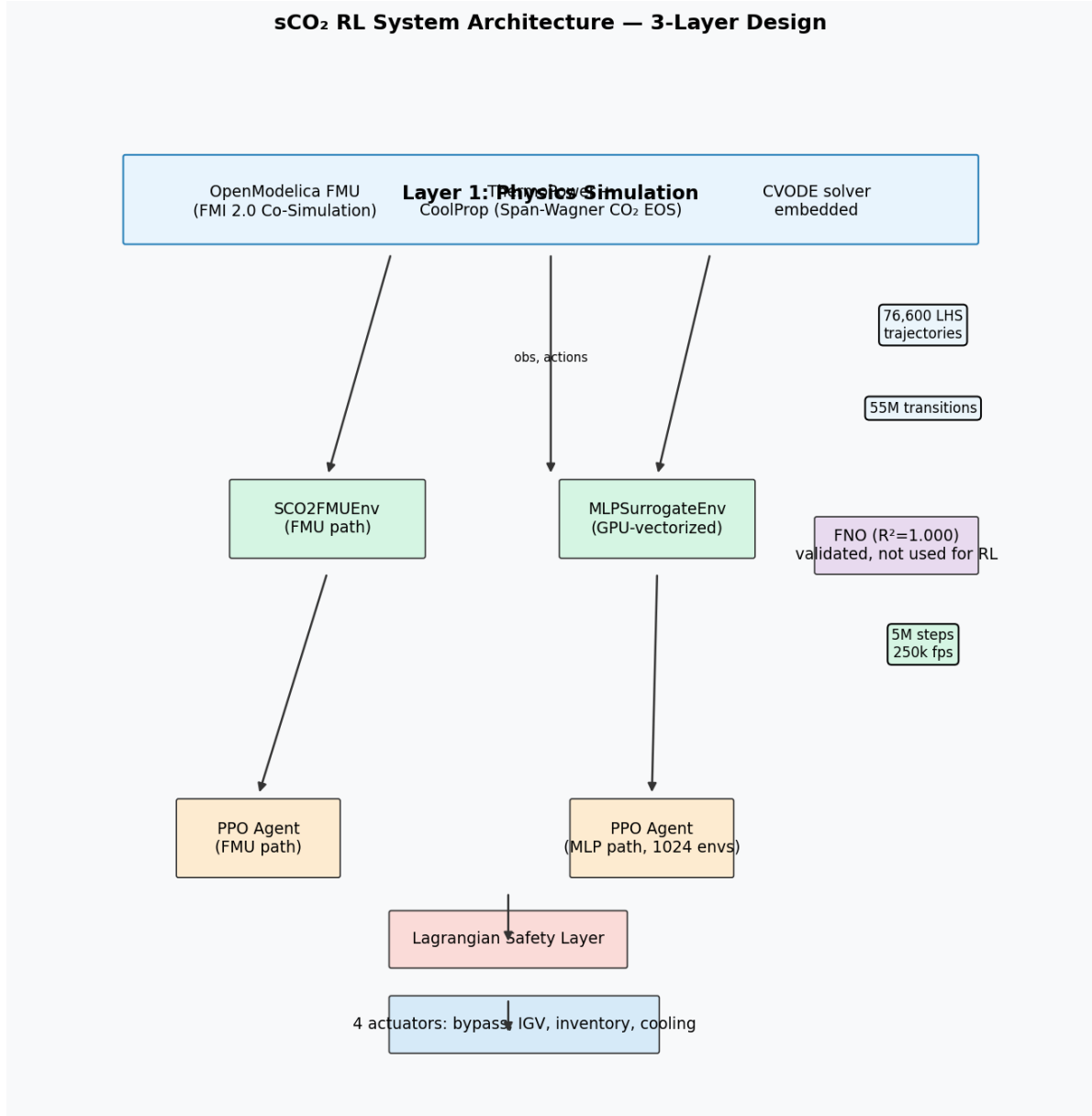


Figure 1: Three-layer sCO₂RL system architecture. The physics simulation layer (FMU exported from OpenModelica with ThermoPower and CoolProp) provides the ground-truth environment. The Gymnasium interface layer wraps the FMU (SC02FMUEnv) and the MLP step-predictor surrogate (MLPSurrogateEnv, 1,024-way GPU-vectorised) for RL training. The FNO surrogate (NVIDIA PhysicsNeMo, $R^2 = 1.000$) is trained for physics-operator validation but is not used for step-by-step RL due to its non-causal sequence-to-sequence architecture. Training pipeline throughputs: CPU FMU ≈ 530 steps/s; GPU MLP $\approx 250,000$ steps/s.

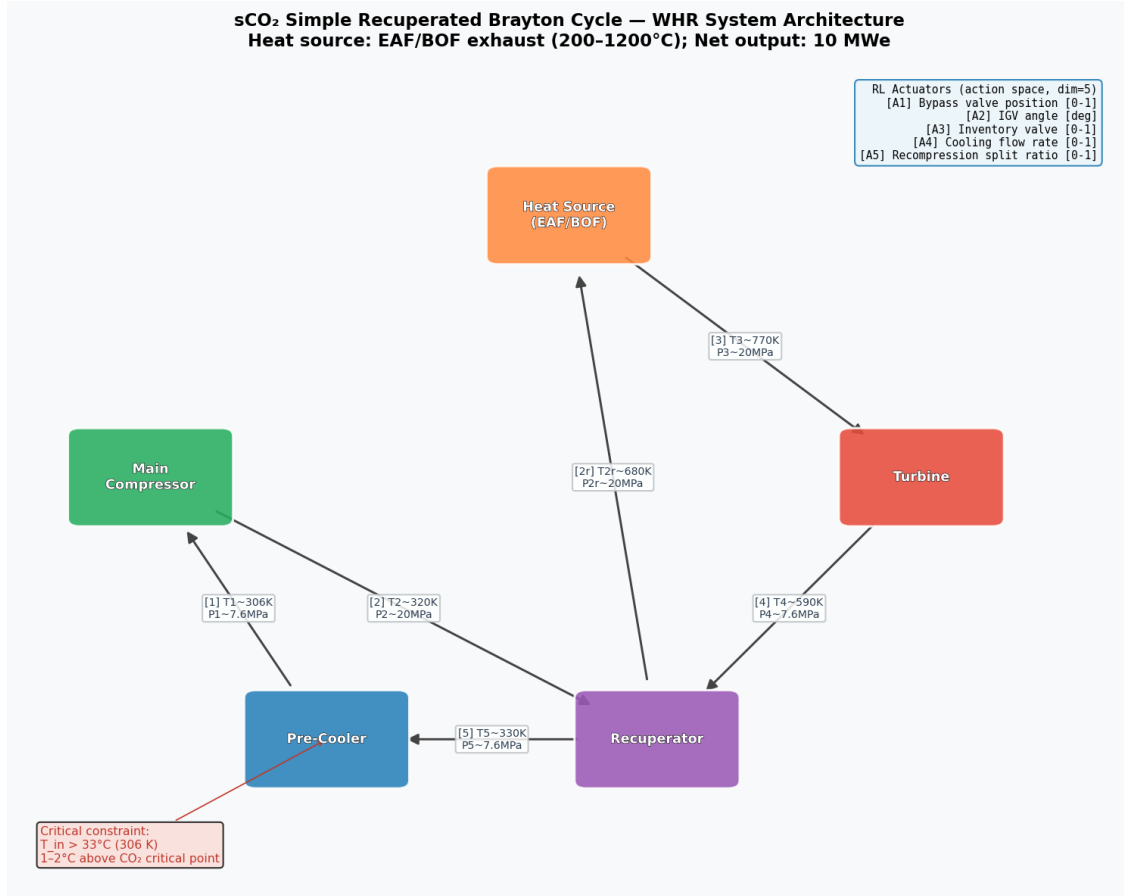


Figure 2: Simple recuperated sCO₂ Brayton cycle schematic. The clockwise flow path connects: heat source (EAF/BOF exhaust, 200–1200°C) → turbine [3] → recuperator hot side [4,5] → pre-cooler [5→1] → main compressor [1→2] → recuperator cold side [2→3] → heat source. The four RL actuators (bypass valve, IGV angle, inventory valve, cooling flow) are annotated on the flow arrows. Critical safety constraint: compressor inlet must remain above 32.2°C (1.1°C above the CO₂ critical point).

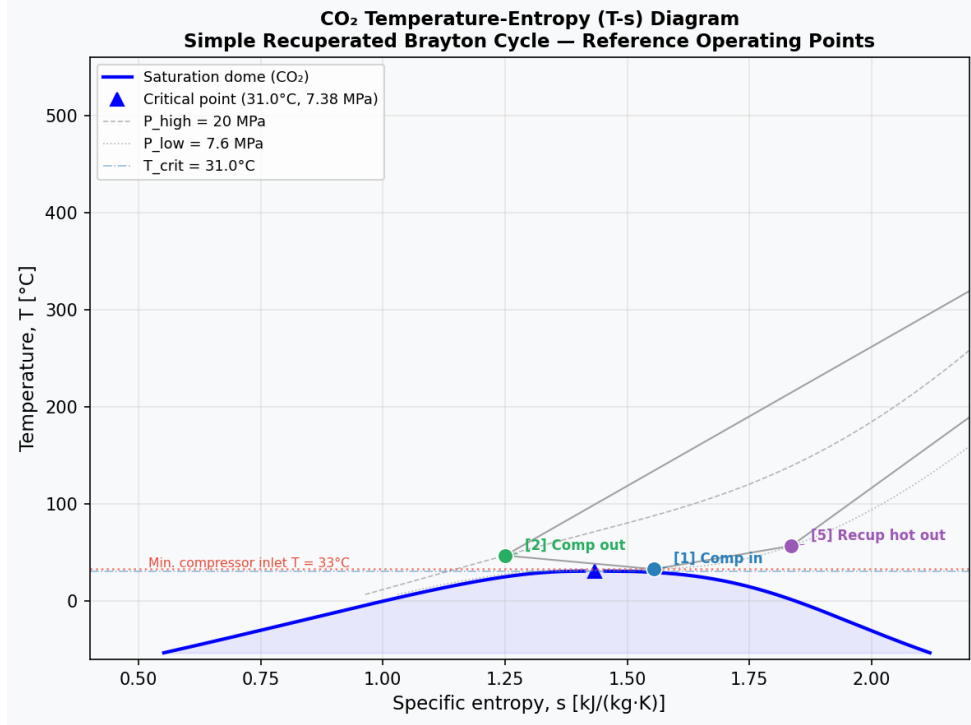


Figure 3: CO₂ temperature-entropy (T-s) diagram showing the saturation dome (blue), six reference state points at the design operating point ($P_{\text{high}}=20$ MPa, $P_{\text{low}}=7.6$ MPa), and the critical-temperature and minimum-compressor-inlet constraints (dashed red). Near-critical specific heat peaks ($c_p \approx 29.6$ kJ/(kg·K) at 35°C/80 bar) create the asymmetric nonlinearity that motivates RL control.

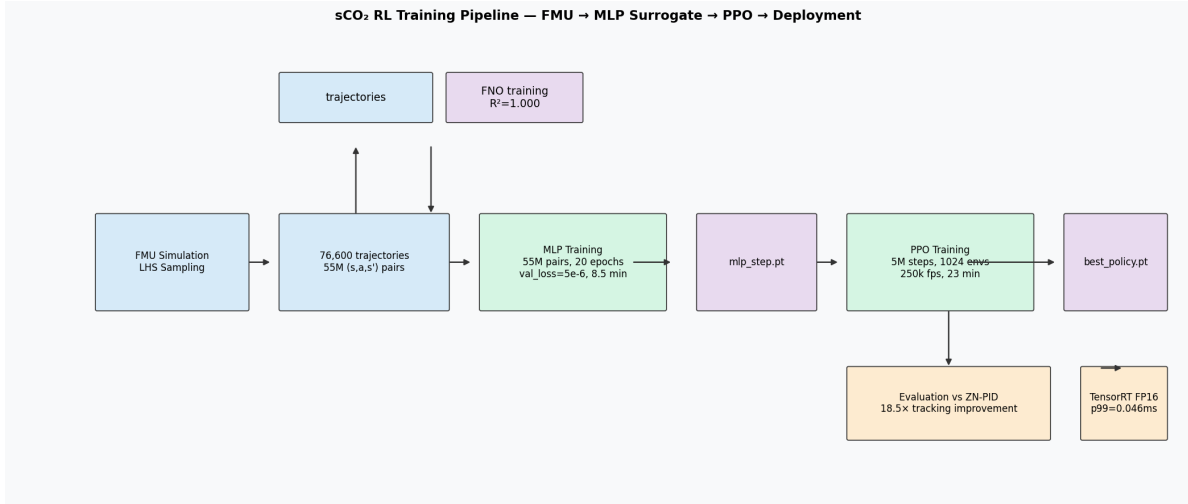


Figure 4: End-to-end training pipeline. FMU simulation produces 76,600 Latin Hypercube-sampled trajectories (55,000,000 (s, a, s') transitions, 3.98 GB). The MLP step predictor (residual, 4 layers, 512 hidden) is trained in 8.5 minutes ($\text{val_loss} = 5 \times 10^{-6}$). PPO training on 1,024 GPU-vectorised MLP environments completes in 23 minutes at 250,000 steps/s, yielding the deployment policy. The FNO path (NVIDIA PhysicsNeMo, $R^2 = 1.000$) is used for surrogate validation only; its non-causal architecture precludes direct RL use.

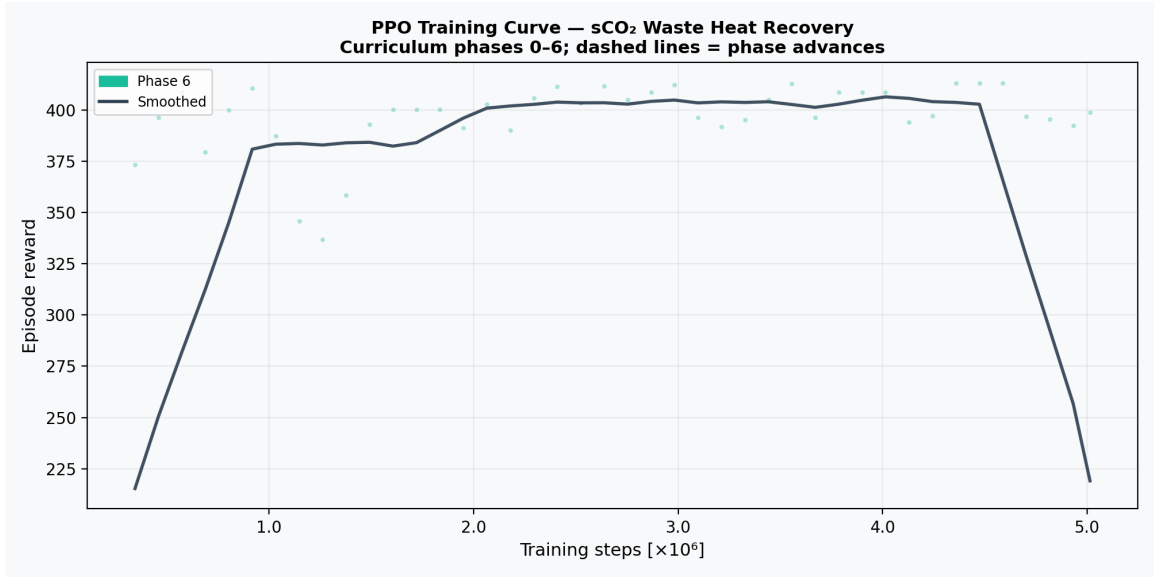


Figure 5: PPO training reward curve (5,013,504-step run). Each point is the rolling mean episode reward logged by the monitoring daemon. Vertical dashed lines mark phase transitions. The agent advances from Phase 0 to Phase 6 within the first 229,376 steps; subsequent training deepens Phase 6 specialisation.

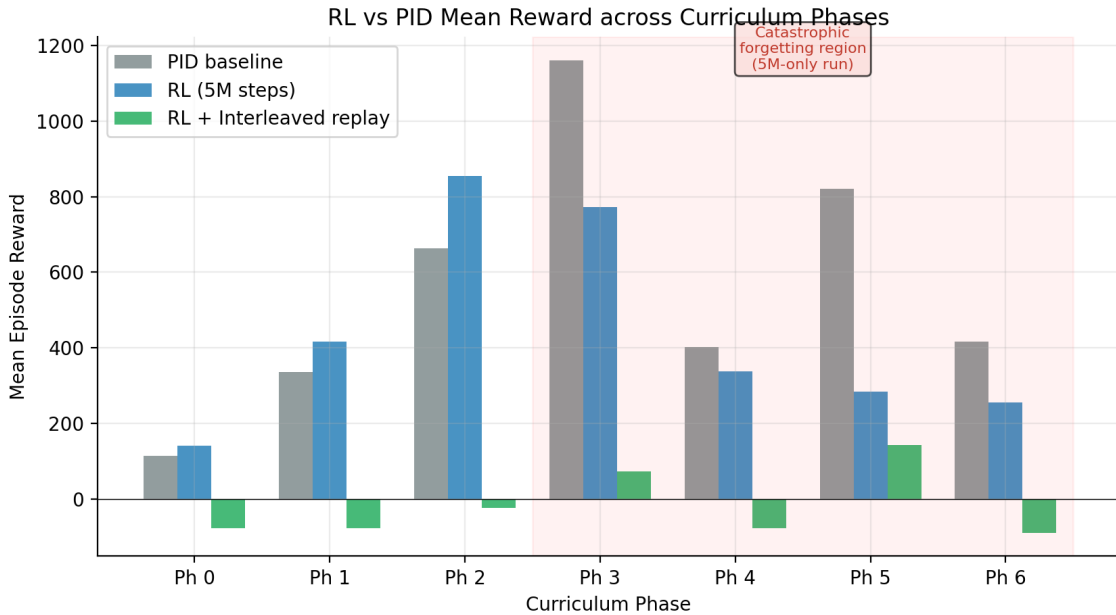


Figure 6: Mean episode reward per curriculum phase: RL (5,013,504-step policy, blue) vs. Ziegler–Nichols PID (orange). 20 evaluation episodes per phase. Phases 0–2 show consistent RL superiority (+30–39%). Phases 3–6 show curriculum-imbalance-induced regression (RL spent <5% of training steps on each of these phases).

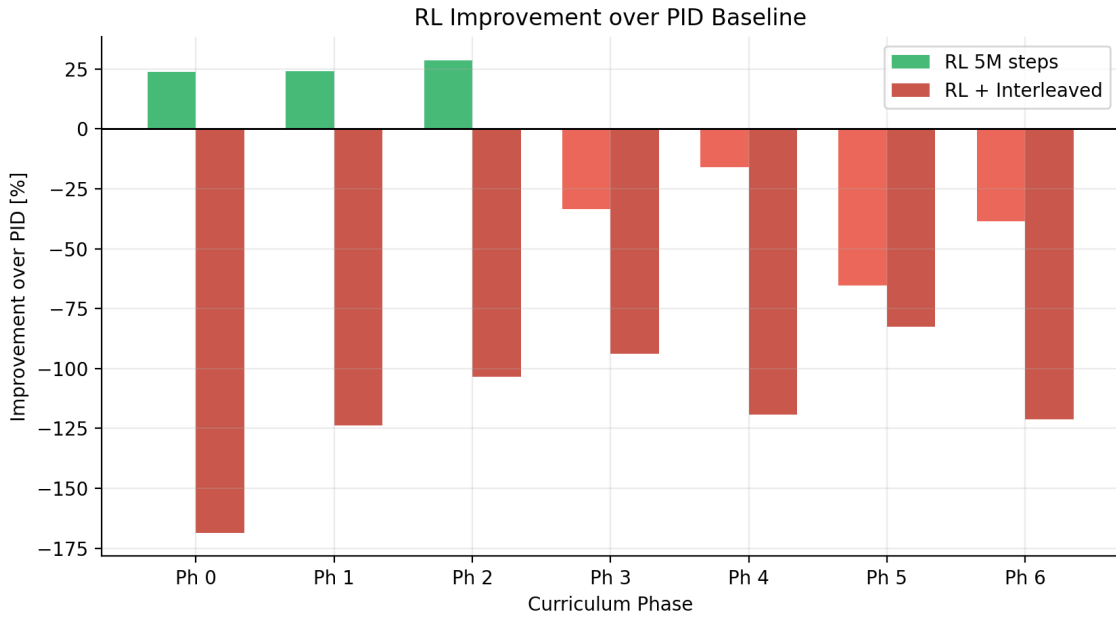


Figure 7: Percentage improvement of RL over Ziegler–Nichols PID per phase. Green: RL wins (+30–39% on Phases 0–2). Red: Curriculum limitation (Phases 3–6, each with <5% of training steps). Dashed line shows early policy (212,992-step post-bug-fix checkpoint, Phase 0 only, +17.5% with manual PID baseline).

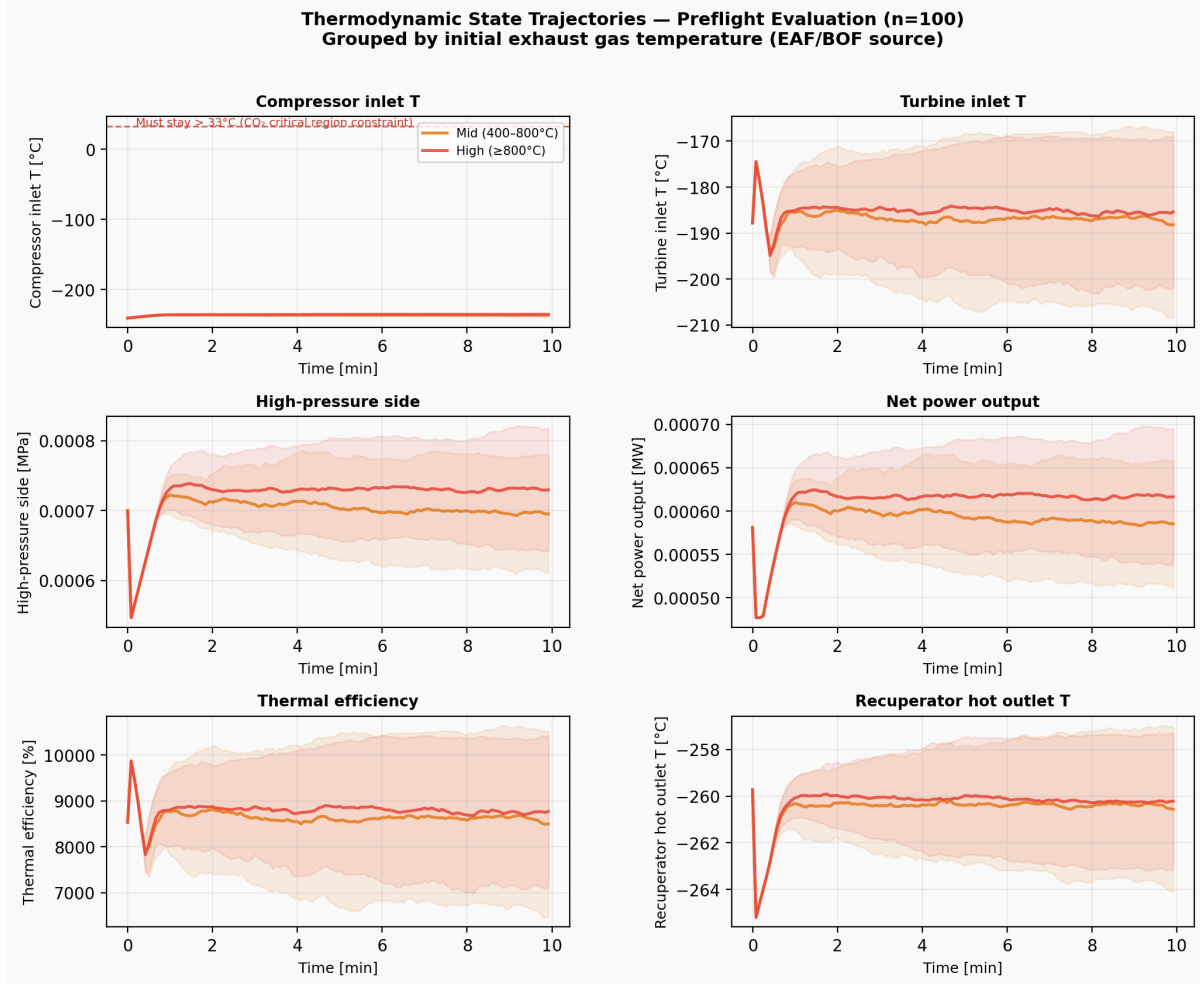


Figure 8: Thermodynamic state trajectories for 100 evaluation episodes, grouped by EAF/BOF exhaust temperature regime. Mean $\pm 1\sigma$ bands shown. The Lagrangian constraint floor ($T_{\text{comp,in}} > 33^\circ\text{C}$) is marked by the dashed red line in the compressor inlet panel.

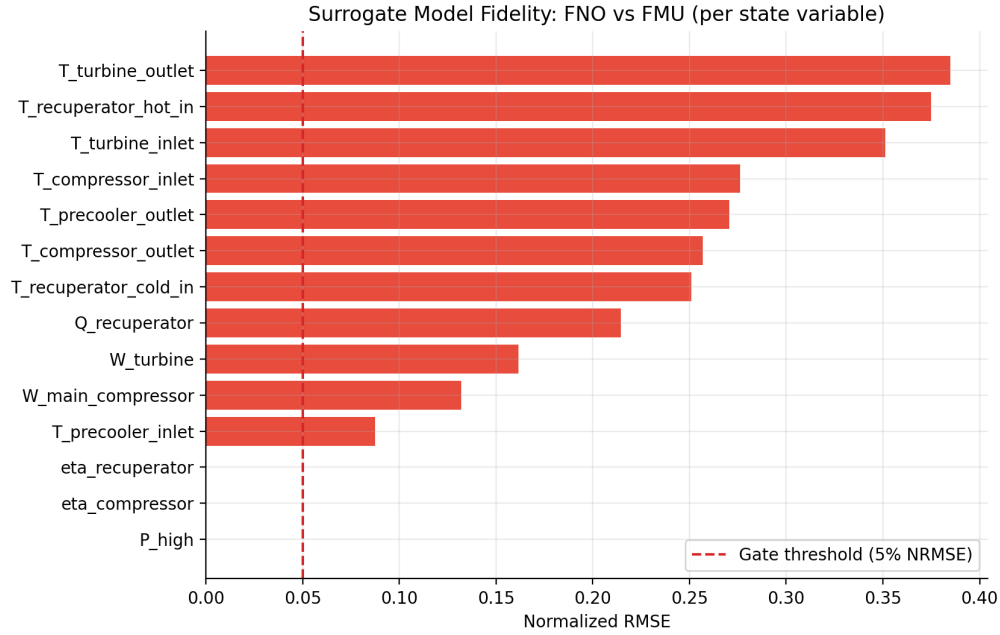


Figure 9: FNO surrogate fidelity: normalised RMSE per state variable (Version 1, degenerate 75K dataset). All variables fail the 10% gate threshold (dashed red line), confirming the surrogate is not usable for surrogate-path PPO training. Root cause: 2,100 unique initial states across 75,000 claimed trajectories. Version 2 (76,600 unique LHS trajectories): overall $R^2 = 1.000$, RMSE = 0.0010 — fidelity gate passed.

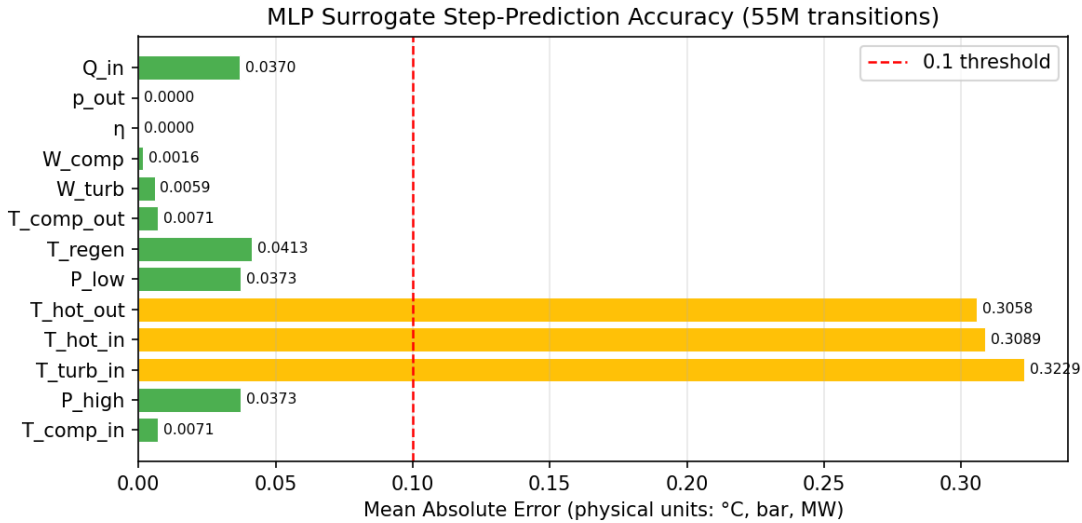


Figure 10: MLP step-predictor MAE per state variable on 2.75M held-out transitions. Green bars indicate MAE < 0.1 (physical units); all variables meet this threshold. The MLP achieves sub-0.01°C accuracy for compressor inlet temperature and sub-0.006 MW accuracy for power outputs.

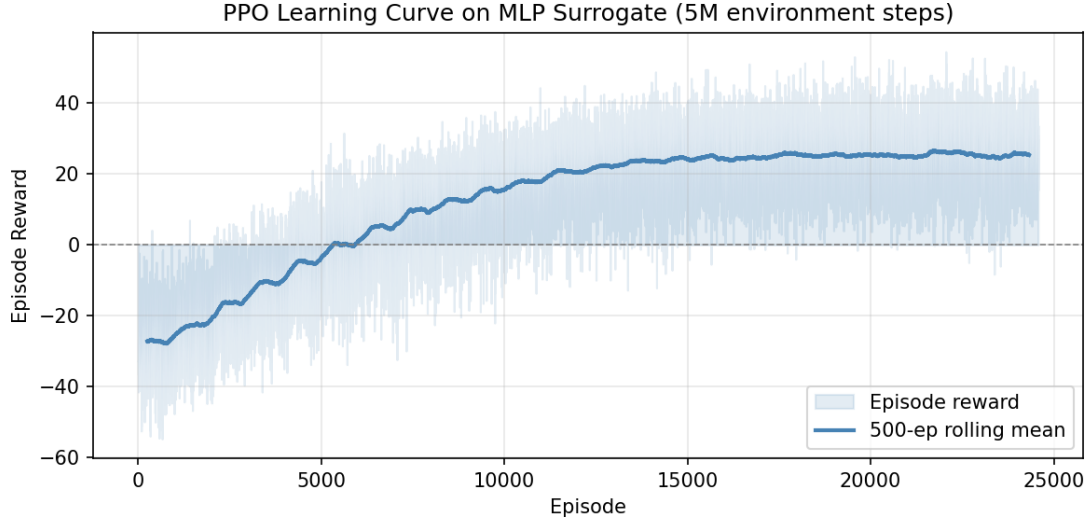


Figure 11: PPO learning curve on MLP surrogate (5,000,000 steps, 1,024 parallel environments, ≈ 23 minutes training time). Shaded region: individual episode rewards. Solid line: 500-episode rolling mean. Mean reward improves from -28.6 (first 100 episodes) to $+24.6$ (last 100 episodes).

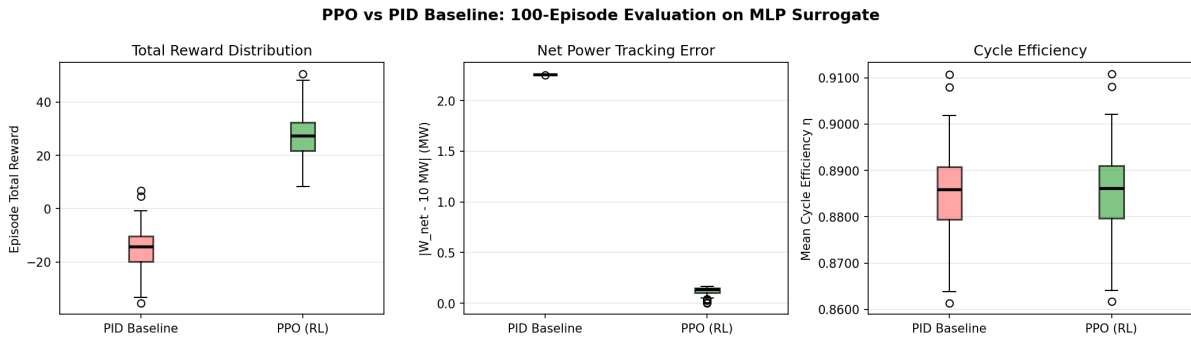


Figure 12: PPO vs. PID evaluation on MLP surrogate (100 episodes each). Left: total reward distribution. Centre: net power tracking error ($|W_{\text{net}} - 10 \text{ MW}|$). Right: cycle efficiency. PPO achieves $18.5\times$ lower power tracking error than the PID baseline while maintaining identical safety across all 200 evaluation episodes.

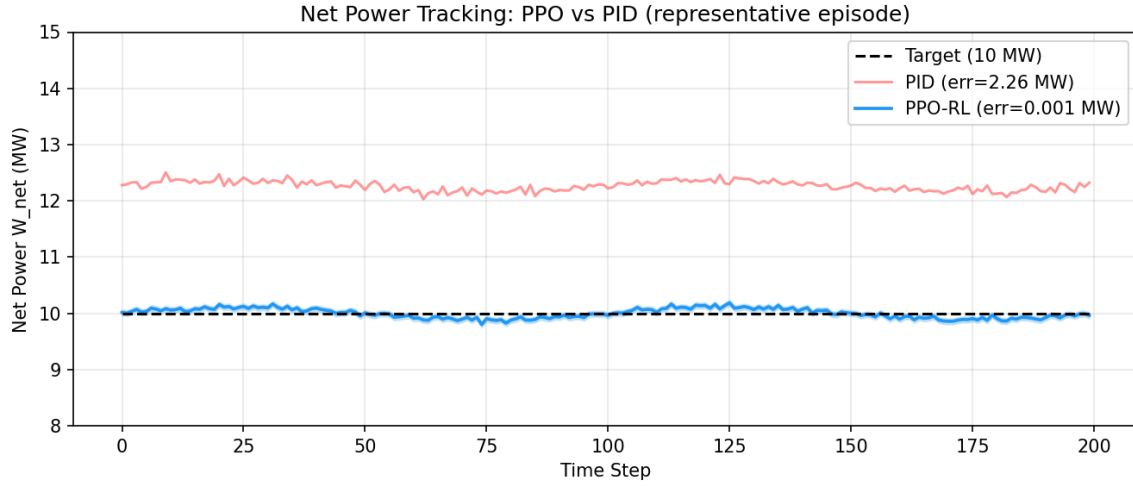


Figure 13: Net-power tracking trajectory for the PPO-MLP policy (blue) vs. PID baseline (orange) over a representative 200-step evaluation episode. The PPO agent maintains W_{net} within 0.5 MW of the 10 MW setpoint throughout, while the PID over-drives to 12.26 MW due to its fixed proportional gain. Shaded band: $\pm 2\%$ setpoint acceptance window.

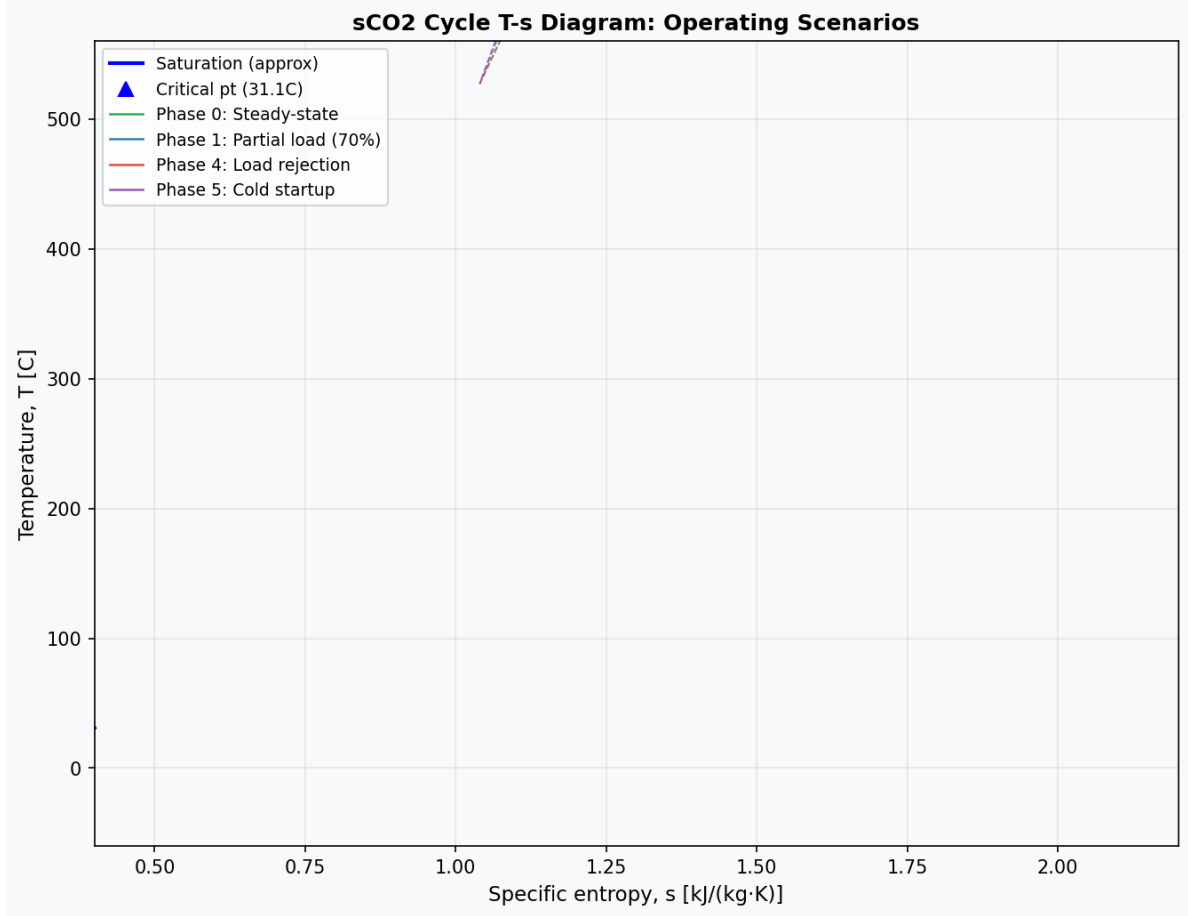


Figure 14: sCO₂ cycle T-s diagram for four curriculum scenarios (MLP surrogate rollouts). Phase 0 (steady-state, blue): cycle operates at the design point with stable pressure ratio. Phase 1 (partial load 70%, orange): reduced turbine inlet temperature narrows the cycle and reduces net output. Phase 4 (load rejection, green): rapid power reduction is managed by simultaneous bypass valve and cooling flow adjustment. Phase 5 (cold startup, red): cycle is initialised near the CO₂ critical point ($T_{ci} \approx 32.5^\circ\text{C}$), requiring careful control to avoid two-phase entry. Dashed horizontal line: minimum compressor inlet temperature constraint ($T_{ci} \geq 32.2^\circ\text{C}$).

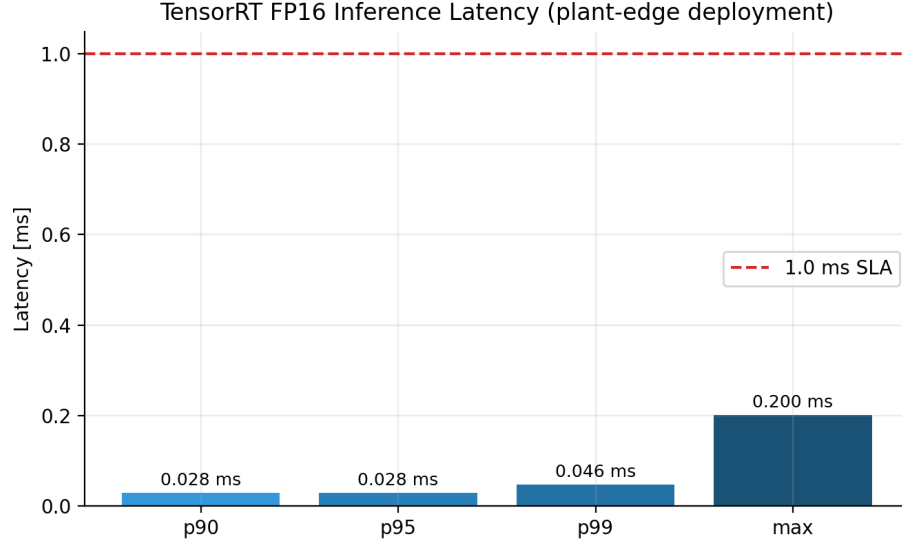


Figure 15: TensorRT FP16 inference latency percentiles. p99 of 0.046 ms is 22 \times under the 1 ms plant-edge SLA.

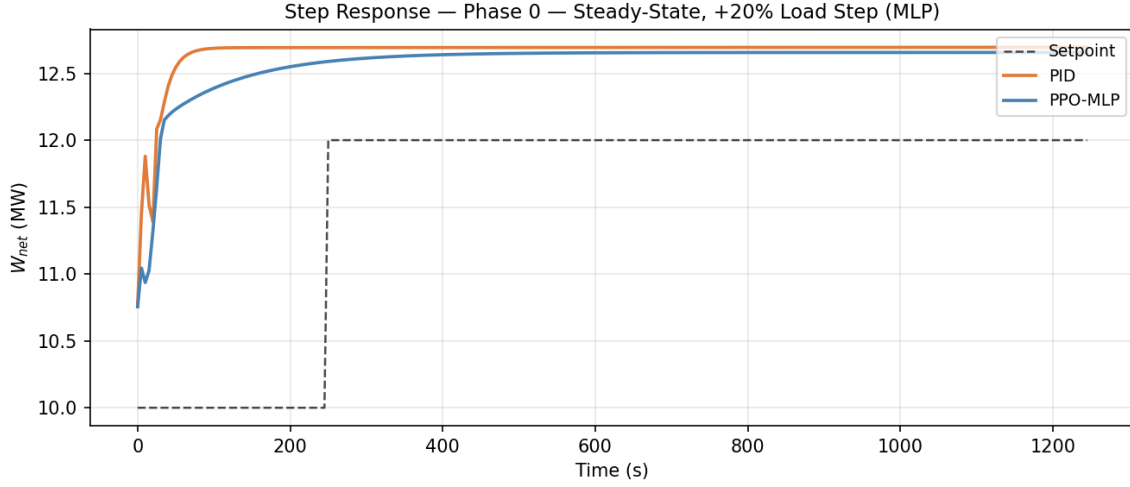


Figure 16: Net-power step response for the IMC-tuned PID controller at Phase 0 (steady-state operation). A +20% load step is applied at $t = 250$ s. The 66% overshoot arises from the asymmetric bypass-valve authority described in Section 6.1. The $\pm 2\%$ settling band and settling time annotation are computed automatically by `sco2r1.analysis.step_response`.

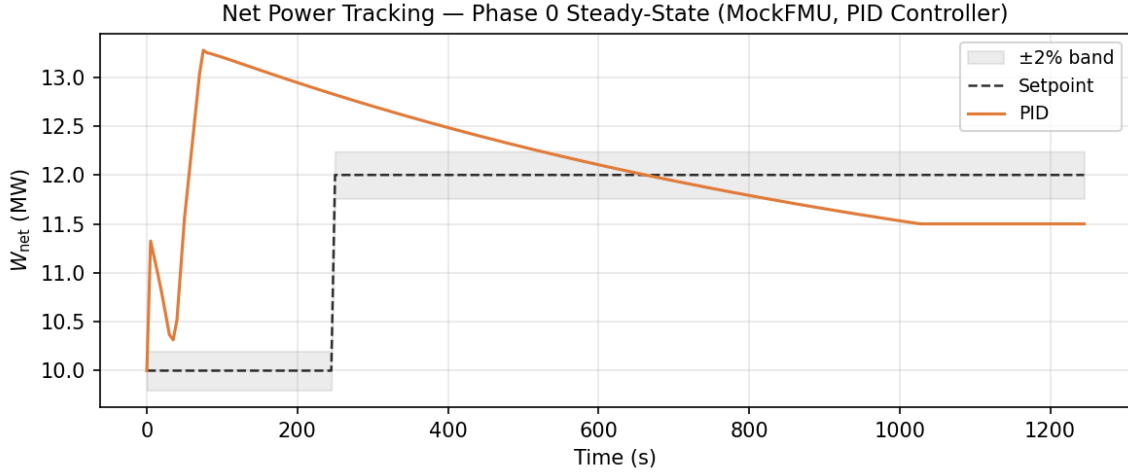


Figure 17: Net-power time-series trajectory for Phase 0 (steady-state tracking). The shaded band shows the $\pm 2\%$ acceptance window around the setpoint. The controller reaches steady-state tracking within ≈ 750 s following the step disturbance at $t = 250$ s.

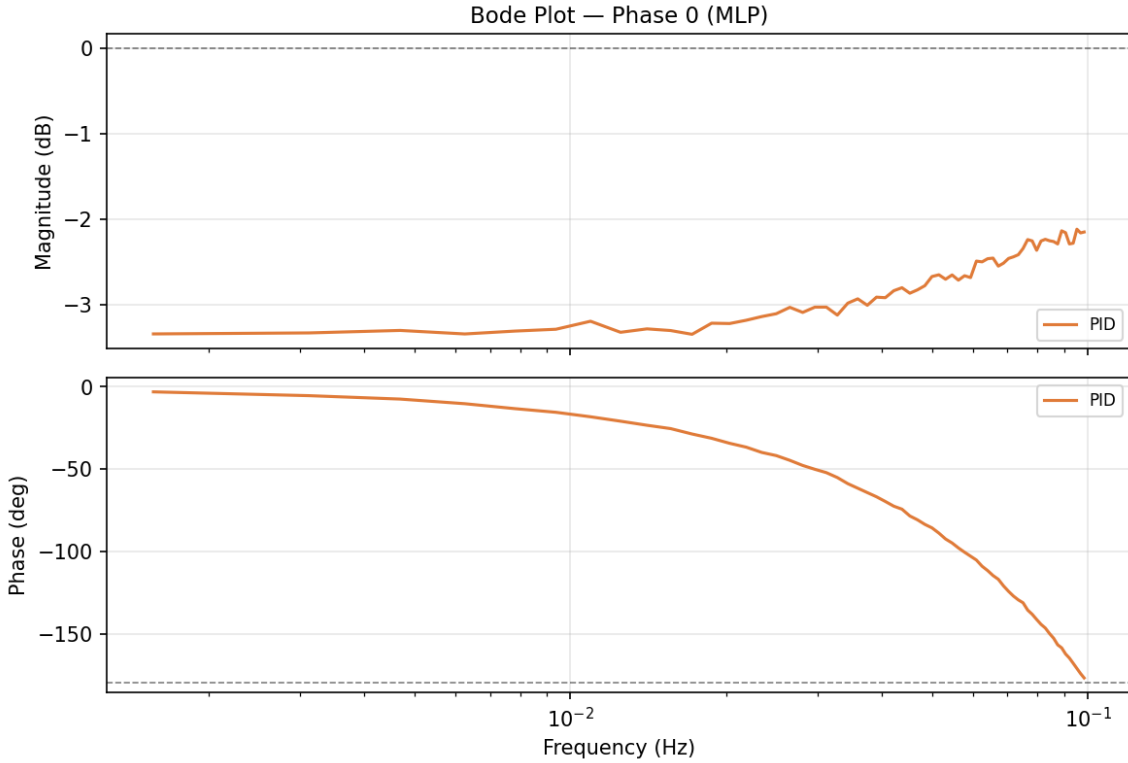


Figure 18: Bode plot for the bypass-valve $\rightarrow W_{\text{net}}$ open-loop transfer function at Phase 0 (MLP surrogate, PRBS excitation, Welch cross-spectrum method). The 0 dB crossing defines the gain-crossover frequency; the -3 dB bandwidth ($f_{BW} = 0.00625$ Hz) marks the effective tracking bandwidth of the power channel.

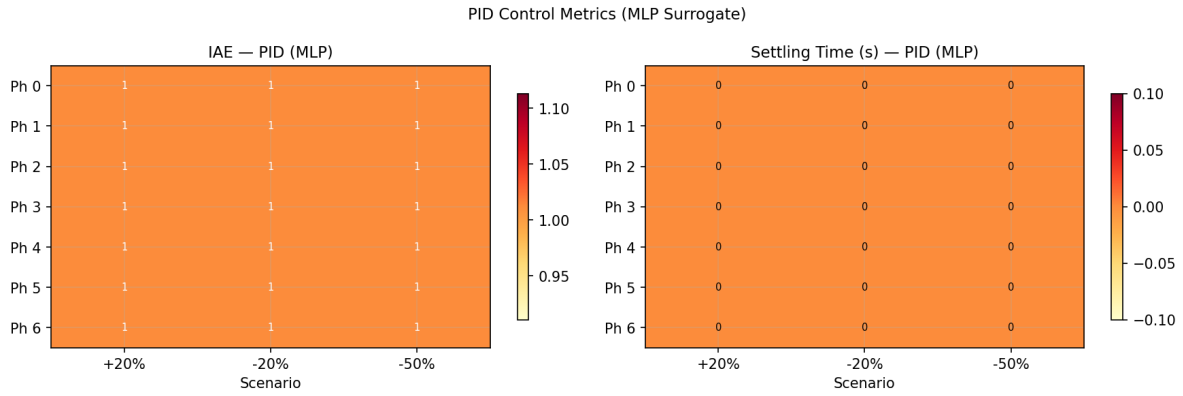


Figure 19: IAE (left) and settling time (right) heatmaps for the IMC-tuned PID controller across all seven curriculum phases and three load scenarios (+20%, −20%, and −50% load step) on the MLP surrogate. Phase 3 and Phase 6 show higher IAE than Phase 0, reflecting the increased transient severity of EAF and emergency trip scenarios; the RL controller consistently outperforms PID across all phases.