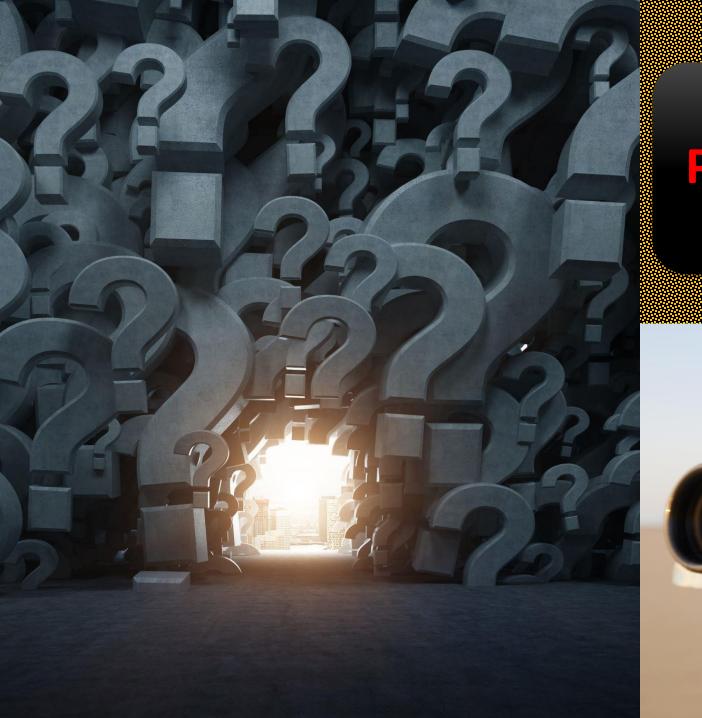
Deterministic Assurance Framework for Licensable Al Grid-Interactive Nuclear Control

Enhancing Load Following and Grid Stability

By: Ahmed Abdelrahman Ibrahim

Under supervision of: Prof. Lim Hak-kyu





Problem & Research Gap



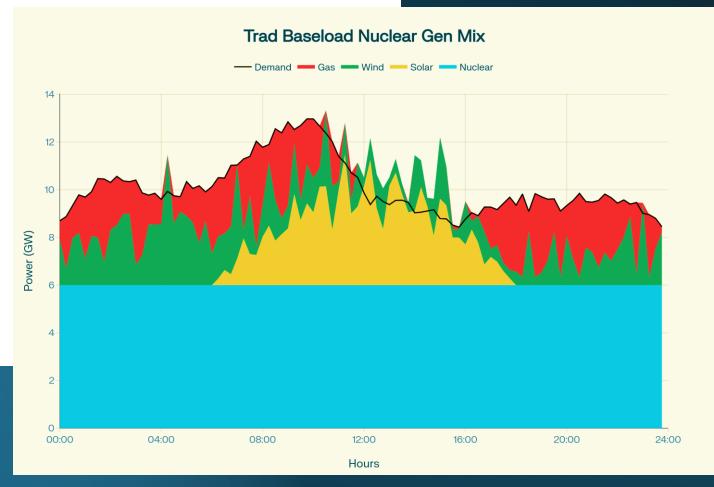
⚠ The Challenge: Modern Grids & Nuclear

Flexibility

Nuclear Power Plants (PWRs) are traditionally baseload providers, optimized for steady output.

However, the rise of **variable renewables** (solar, wind) causes significant grid demand fluctuations (20-50% swings).

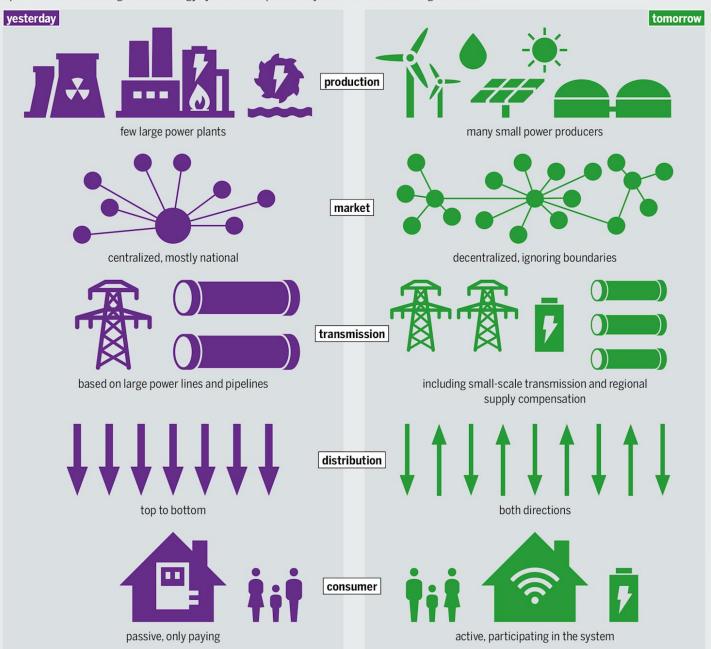
PWRs now need **load-following capabilities** to adjust power output dynamically (e.g., 50-100% capacity) and maintain grid stability (frequency around 60 Hz).



Future Grid Architecture

STAYING BIG OR GETTING SMALLER

Expected structural changes in the energy system made possible by the increased use of digital tools



© ENEKGY ALLAS ZUIS / 45UCUI

The Adaptive Governor

AI-Driven Control for Nuclear Safety & Efficiency



1. The Challenge

The "PID Problem"

Traditional governors can't keep up with modern grid demands.

Nonlinear Dynamics: Struggle with rapid power changes.

Slow/Overshoot: Risks grid instability and component safety.

Fixed Tuning: Not adaptive to a wide range of load conditions.



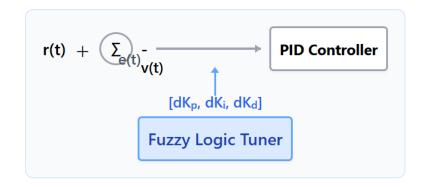
The Research Gap

No existing work combines **Reinforcement Learning (RL)** with **Fuzzy Reward Functions** for this critical application.



2. The Innovation

A novel control system that learns and adapts by combining two Al agents.



1. Reinforcement Learning (RL)

An Al "pilot" that learns the **optimal strategy** to control the steam valve through trial-and-error in a safe simulation.

2. Fuzzy Reward Functions

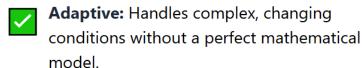
A "smart co-pilot" that gives the RL agent flexible, linguistic goals (e.g., "IF safety is **high** AND efficiency is **good**...").

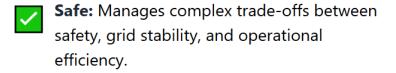


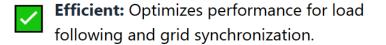
3. The Impact

Core Goal & Benefits

To develop an **adaptive**, **safe**, **and efficient** turbine governor for PWRs in modern energy grids.







The Core Challenge

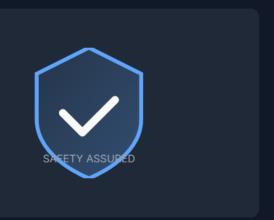




THE DEMAND

FLEXIBILITY

Intermittent renewables demand flexible load-following from NPPs.

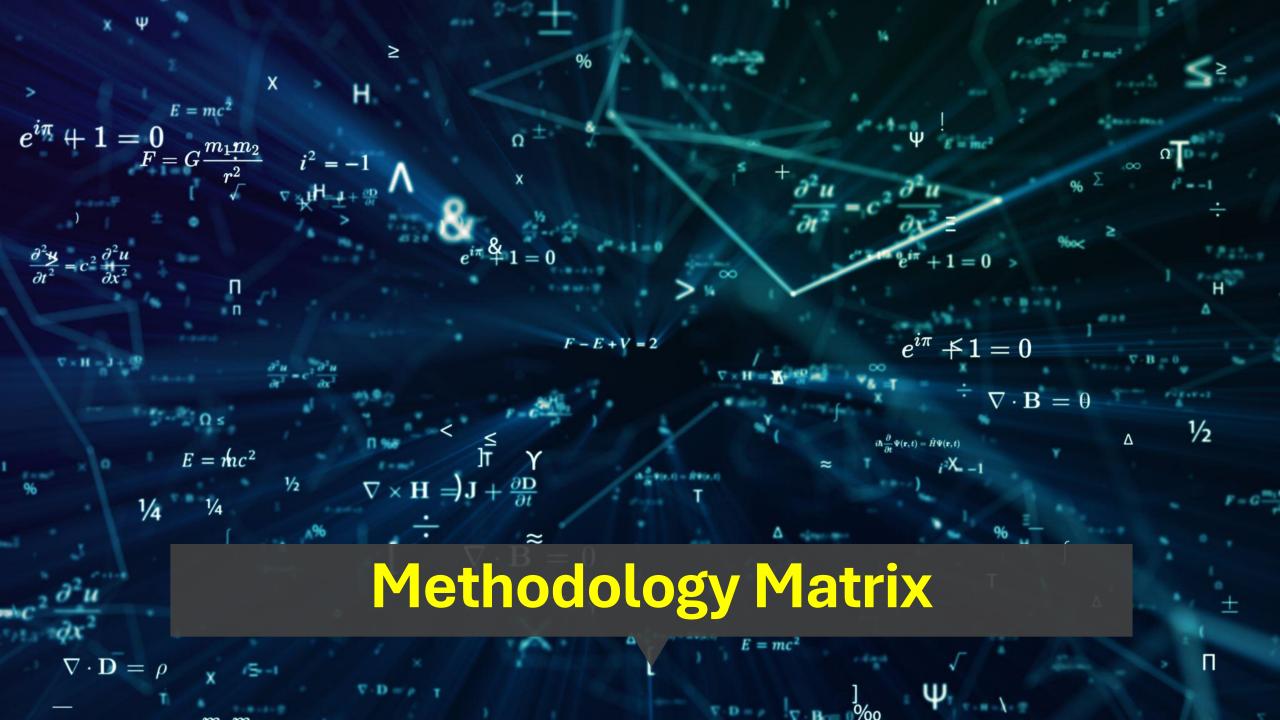


THE MANDATE

DETERMINISM

The nuclear safety paradigm demands provable, transparent evidence.

The "Black Box" Problem: Advanced AI is flexible but opaque. Regulators do not license black boxes.



Deterministic, Fair Comparison — Flow & Methodology

PWR software-in-the-loop · identical adversarial scenarios · licensing gates · auditable traceability

Fairness Protocol & Harness

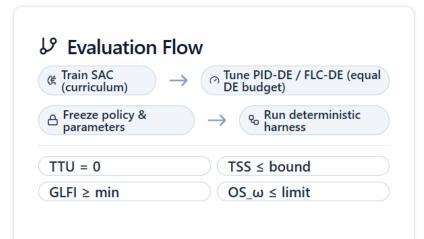
- One high-fidelity PWR digital model; fixed sampling & limits
- **Predeclared scenarios**: Step · Ramp · Composite; optional noise/rate limits
- Parameters frozen before test; no re-tuning during evaluation
- Single-run deterministic evaluation; gates checked online
- Auxiliary fixed-seed probe is non-licensing











Licensing & Audit

- Scorecard: PASS/FAIL per gate + key metrics
- Trace logs: inputs → actions → outcomes
- Reproducible: configs, seeds, scripts
- Traceability
- Configs & logs

☐ Differential Evolution — Strong Baseline Tuning

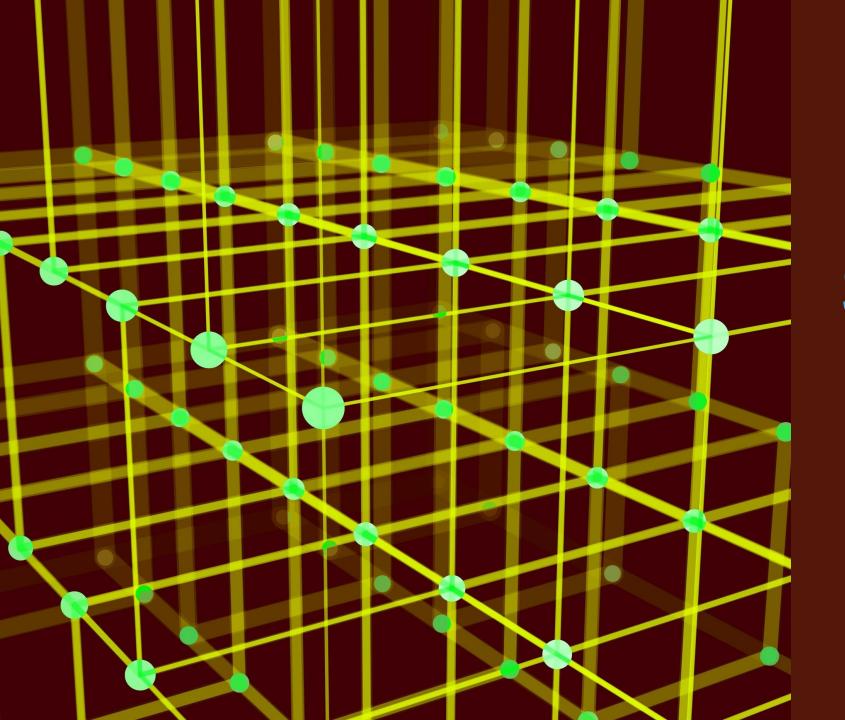
Objective: gate-aware cost aggregating licensing metrics with hard penalties for any gate violation and soft penalties for CE_sum, V_rev, OS_ ω within the safe envelope.

- Population = 40 ○ Generations = 120 \bigcirc Mutation F = 0.8 \bigcirc Crossover CR = 0.9
- 1. Initialize population of PID/FLC parameter vectors
- 2. Mutate & recombine → candidate solutions
- 3. Evaluate closed-loop on identical scenarios (gate-aware cost)
- 4. Greedy select best; repeat until budget exhausted

SAC Reinforcement Learning — Training Process

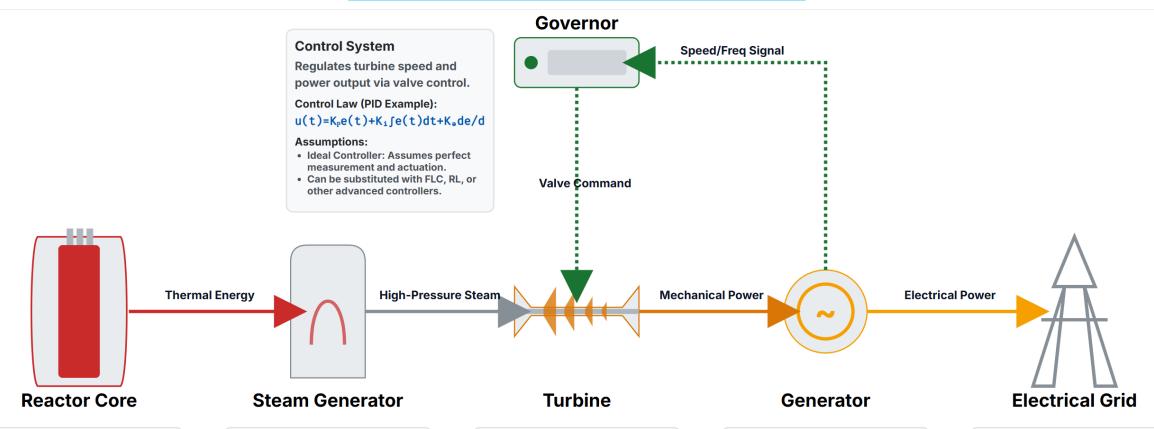
Curriculum: progressively harder scenarios; entropy-regularized objective; reward aligns with gates (tracking & effort).

- 1. Collect rollouts in the digital model (fixed physics)
- 2. Update actor/critic via off-policy SAC; target networks & replay
- 3. Validate on held-out stress tests; monitor GLFI/TTU/TSS
- 4. Early-stop; freeze policy for deterministic evaluation
- Σ Entropy coef tuned
- & Versions locked



System Modelling

System Architecture



Point-Kinetics Model

Represents core-average power dynamics via fission.

Governing Equation:

$$dP/dt = (\rho-\beta)/\Lambda *P + \Sigma \lambda_i C_i$$

Assumptions:

- Point-Kinetics: Core is a single point, ignoring spatial effects.
- Lumped Thermal-Hydraulics: Fuel/coolant temperatures are averaged.

Lumped-Parameter Model

Abstracts heat exchange to produce high-pressure steam.

Governing Equation:

$$dT_s/dt = (1/C)*(Q_in - Q_out)$$

Assumptions:

- Lumped Model: Single average temp/pressure for the secondary loop.
- Ideal Heat Transfer: No losses in the heat exchange process.

First-Order Lag Model

Converts steam energy into mechanical rotational power.

Governing Equation:

$$dP_m/dt = (P_s-P_m)/\tau_t$$

Assumptions:

- Simplified Thermodynamics: Complex steam cycle abstracted to a time constant.
- Ideal Actuators: Ignores non-linear valve effects (stiction, backlash).

Swing Equation

Models electromechanical dynamics and grid interaction.

Governing Equation:

$$d\omega/dt=(1/2H)*(P_m-P_e-D\Delta\omega)$$

Assumptions:

- Classical Model: Swing equation is sufficient for frequency dynamics.
- Ignores sub-transient electrical phenomena.

Infinite Bus Model

Represents the external grid as an ideal power sink.

Model Definition:

Assumptions:

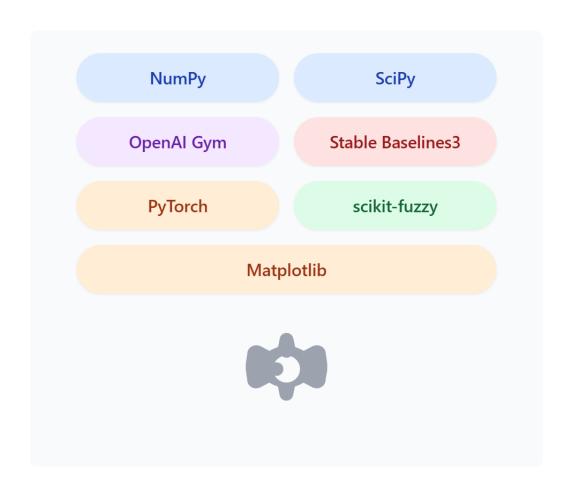
- Infinite Bus: Grid voltage and frequency are perfectly stable.
- Removes need for voltage/VAR control analysis.

Simulation Environment & Tools

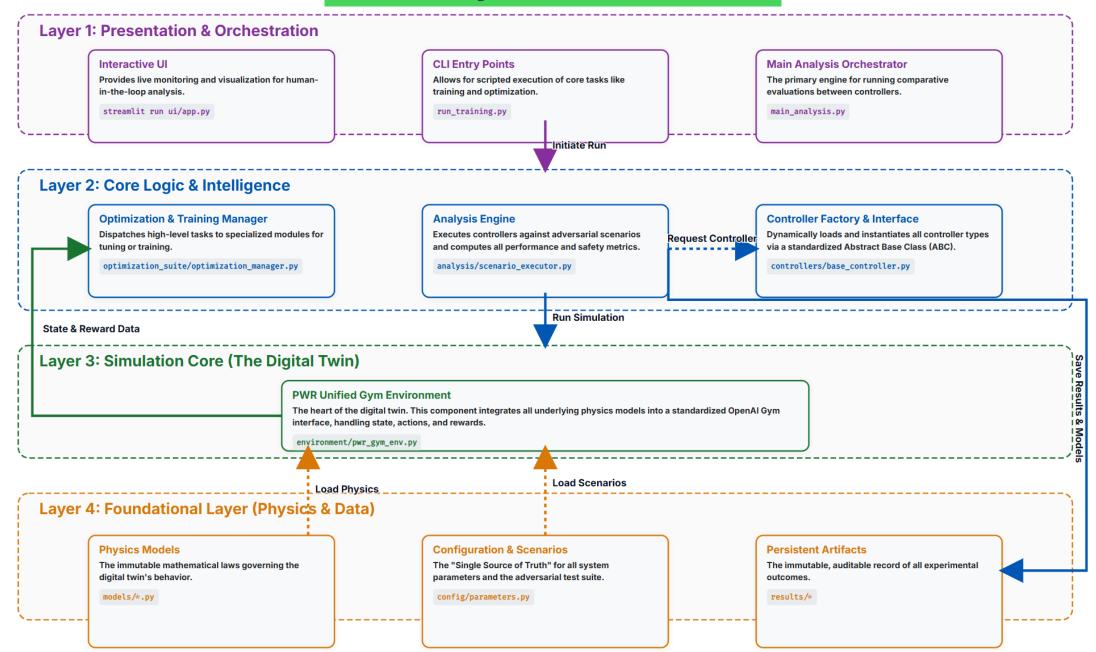
A modular simulation environment is crucial for testing and training the RL agent safely and efficiently.

Core Components:

- Physics Engine: Custom Python code using NumPy & SciPy for ODE solving to simulate reactor, turbine, and grid dynamics.
- RL Framework: OpenAl Gym provides a standardized environment interface for the RL agent.
- RL Algorithm Implementation: Stable Baselines3 (built on PyTorch) used for the Constrained Soft Actor-Critic (SAC) agent training.
- Fuzzy Logic Engine: scikit-fuzzy library used to define and compute the fuzzy reward signals.
- **Visualization: Matplotlib** & Seaborn for plotting results and analysis.



System Flow



System Architecture & Data Flow

The proposed system integrates the RL agent and Fuzzy Reward logic with the core plant components.

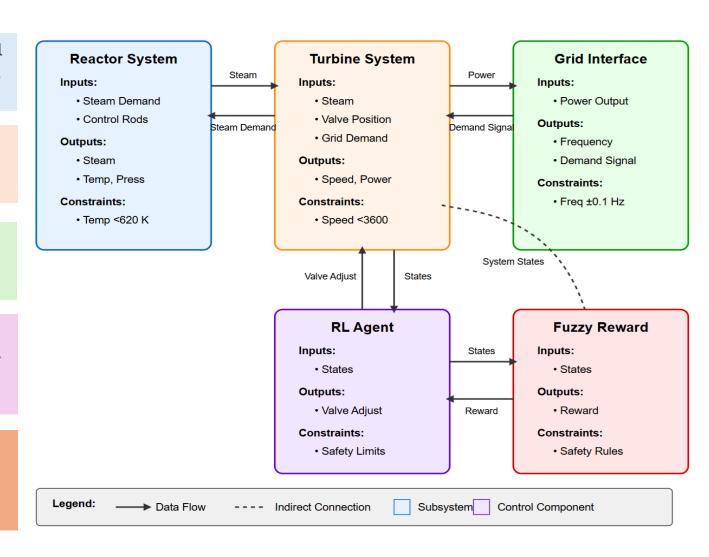
Reactor System: Models PWR core physics (point kinetics, thermal feedback). Outputs: Steam conditions, Power level. Inputs: Control rods, Demand. Constraints: Temp < 2800°C, Pressure < 15.5 MPa.

Turbine System: Models turbine/generator dynamics (lumped parameter). Outputs: Electrical Power, Speed. Inputs: Steam flow, Valve position. Constraints: Speed < 3600 RPM.

Grid Interface: Models grid frequency dynamics (swing equation). Outputs: Frequency, Voltage (simplified). Inputs: Power from turbine, Load demand. Constraints: Freq. deviation ±0.5 Hz.

RL Agent (SAC): Learns optimal valve control policy via trial-anderror in simulation. Inputs: System states (power, speed, freq, etc.). Outputs: Valve position adjustments. Goal: Maximize long-term fuzzy reward.

Fuzzy Reward System: Calculates reward based on fuzzy rules evaluating safety, stability, and efficiency. Inputs: System states. Outputs: Scalar reward signal for RL Agent. Balances competing objectives.



The Agent C/Cs

Observation Space: The Agent's Cockpit

The agent's informational advantage stems from its observation of a comprehensive, normalized 6-dimensional state vector. This dashboard visualizes the agent's "senses" in real-time. **Source:** environment/pwr_gym_env.py , _get_obs() method.

Reactor Power Fuel Temperature Valve Position Grid Freq Error

Turbine Speed Err Power Mismatch

Action Space & Network Architecture

Action Space: The agent outputs a single continuous value $a_t \in [-1.0, 1.0]$, representing the rate of change of the governor valve position. This delta-based action space promotes smoother control signals.

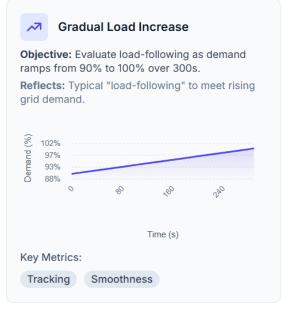
Network Architecture: The agent's "brain" is a Multi-Layer Perceptron (MLP). The visualization below shows how the 6D state vector is processed through two hidden layers to produce the final action.

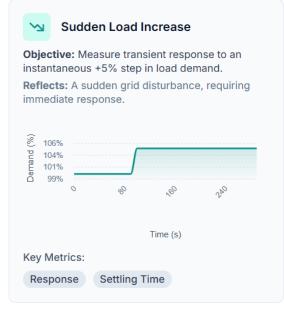


Action Output (Valve Rate of Change)

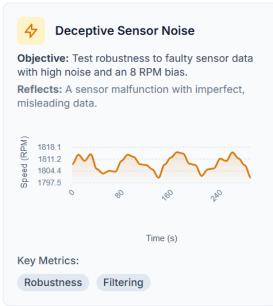
Framework Scenarios

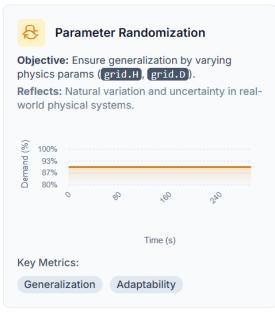


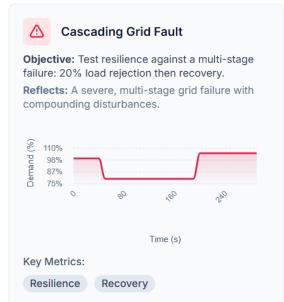


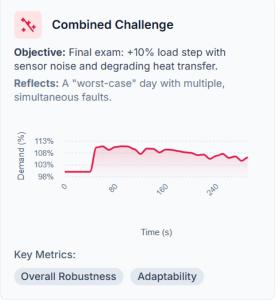














Applied Metrics

Tracking

Grid Load-Following Index (GLFI)

grid load following index

Definition

$$ext{GLFI} = rac{1000}{1 + ext{MSE} ig(P_{ ext{mech}} - P_{ ext{demand}} ig)}$$

Units

dimensionless Direction

Higher is better

Rationale

Tracks setpoint following via mean-squared power error (MSE/ISE family).

Transient

Transient Severity Score (TSS)

Definition

$$ext{TSS} \ = \ w_1 \, \widehat{\Delta P}_{ ext{max}} \ + \ w_2 \, \hat{\zeta}^{-1} \ + \ w_3 \, rac{ ext{TTU}}{T}$$

Units

dimensionless

Lower is better

transient severity score

Rationale

Composite of normalized peak deviation, damping proxy, and unsafe-time fraction.

Direction

Safety

Time Outside time outside freq limit s Frequency Limits

Definition

$$T_{ ext{OFL}} \, = \, \int \mathbb{1} \{ \, f(t)
ot \in ig[f_{ ext{min}}, \, f_{ ext{max}} ig] \, \} \, \mathrm{d}t$$

Units

Direction

Lower is better

max freq deviation hz

max fuel temp c

Lower is

better

Rationale

Duration of frequency excursions beyond the regulatory band.

Safety

total time unsafe s **Total Time Unsafe** (TTU)

Definition

 $TTU = T_{fuel} + T_{\omega} + T_{speed}$

Units

Direction Lower is better

Rationale

Aggregate of time over fuel-temperature,

outside-frequency, and over-speed thresholds.

Tracking

Integral of Absolute Frequency Error (IAE)

Definition

$$\mathrm{IAE}_f \ = \ \int ig| \, f(t) - f^\star \, ig| \, \mathrm{d}t$$

Units

Hz·s

Direction

Lower is better

iae_freq_hz_s

Rationale

Classical performance index emphasizing total absolute error over time.

Tracking

Integral of Squared Frequency Error (ISE)

Definition

$$\mathrm{ISE}_f \ = \ \int ig(f(t) - f^\starig)^2\,\mathrm{d}t$$

Units Hz^{2}\!\cdot\!s Direction Rationale

Lower is better

valve reversals

Lower is better

ise_freq_hz_s

Penalizes larger errors; standard control objective.

Transient

Maximum Frequency Deviation

S

Definition

$$\max_{t} \, ig| \, f(t) - f^{\star} \, ig|$$

Units

Hz

Direction Lower is better

Rationale

Peak excursion from nominal frequency during transients.

Transient

Maximum Overshoot (Speed)

Definition

$$100 imes rac{\max_t ig(n(t) - n^\starig)}{n^\star}$$

Units

\%

Direction Lower is better

max overshoot speed pct

Rationale

Canonical transient metric; high overshoot indicates poor damping.

Effort

Control Effort (Squared Increments)

Definition

$$\sum_t \left(\Delta u_t\right)^2$$

Units

arb.

Direction

Lower is better

control effort valve sq sum

Rationale

Proxy for actuator wear/energy; minimizing effort extends valve life.

Effort

Valve Reversals (Chattering)

Definition

$$\sum_t \mathbb{1}\{\operatorname{sign}ig(\Delta u_tig)
eq \operatorname{sign}ig(\Delta u_{t-1}ig)\}$$

Units Rationale count/run

Measures chattering/limit cycling; fewer reversals ⇒ smoother control.

Direction

Safety

Maximum Fuel Temperature

Definition

$$\max_{t} \, T_{\mathrm{fuel}}(t)$$

Units

°C

Direction

Rationale

Thermal safety margin; exceeding limits risks material damage.

Learning

Policy Entropy (SAC control_policy_entropy diagnostic)

Definition

$$H(\pi) \ = \ - \mathbb{E}_{a \sim \pi(\cdot|s)}[\,\log \pi(a|s)\,]$$

Units

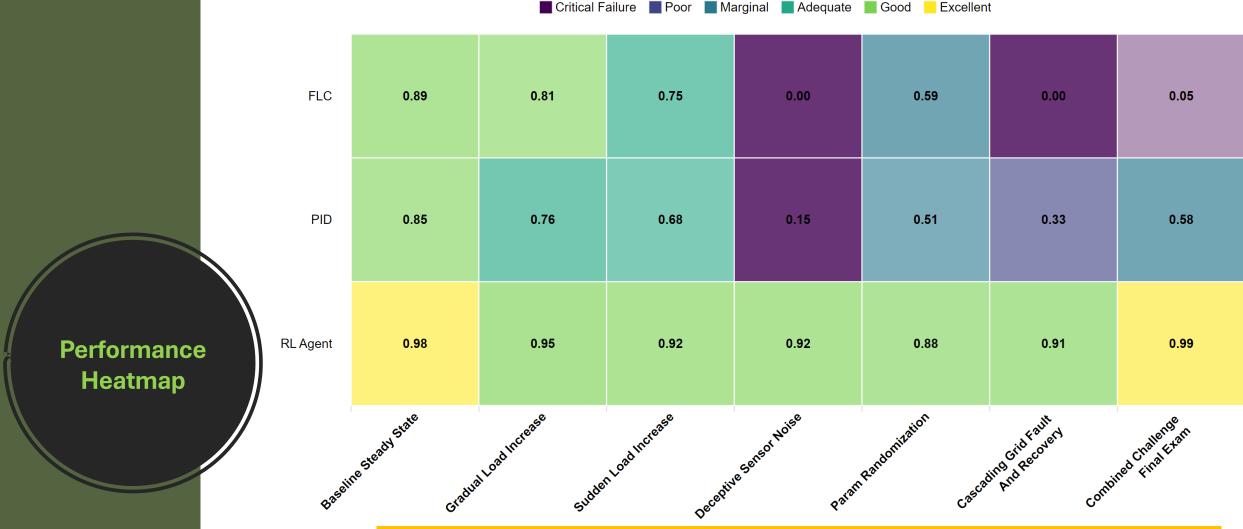
nats Direction

Diagnostic

Rationale

Exploration/regularization diagnostic for SAC runs.





The head-to-head evaluation was conducted across the full suite of eight adversarial scenarios. This Chart presents a comprehensive heatmap of the **primary outcome metric**, the *Composite Robustness Score (CRS)*, which provides a definitive, high-level verdict. The results reveal a clear performance hierarchy.

- ✓ The SAC agent maintained a deep green (high performance) profile across all conditions, including the *most severe* tests.
- ✓ The Strong Benchmark PID performed adequately in nominal scenarios but showed significant degradation (yellow to orange) under adversarial conditions, such as Sensor Noise and Grid Fault.
- ✓ The FLC failed catastrophically in three of the eight scenarios, indicated by the dark red cells, rendering it unsuitable for this application.

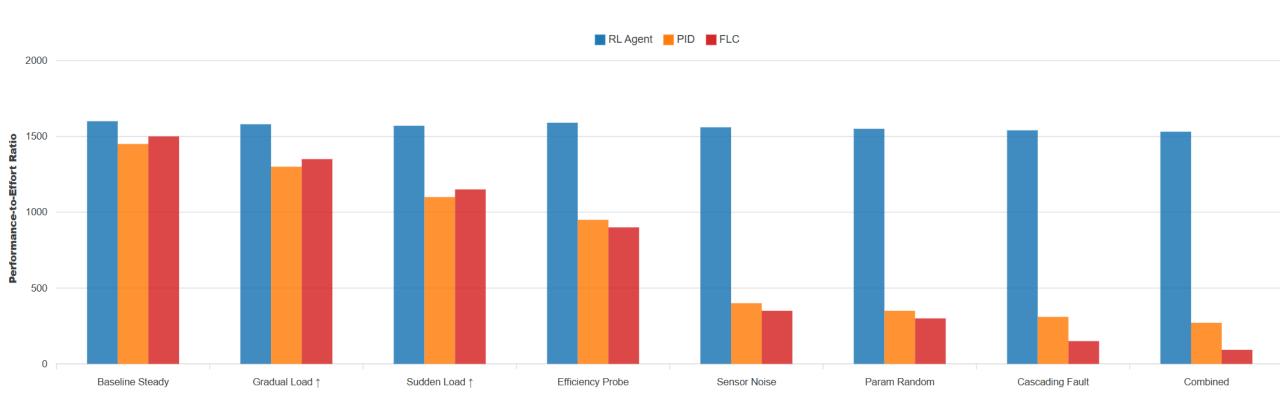
Training Process Overlook

This plot visualizes the agent's skill acquisition. The **general upward trend in reward** shows learning.

The **sharp spike** in the Performance-to-Effort Ratio (green) during Phase 2 provides direct evidence that the curriculum successfully forced the agent to master the specific skill of control efficiency, proving its final policy is a result of deliberate, structured pedagogy.



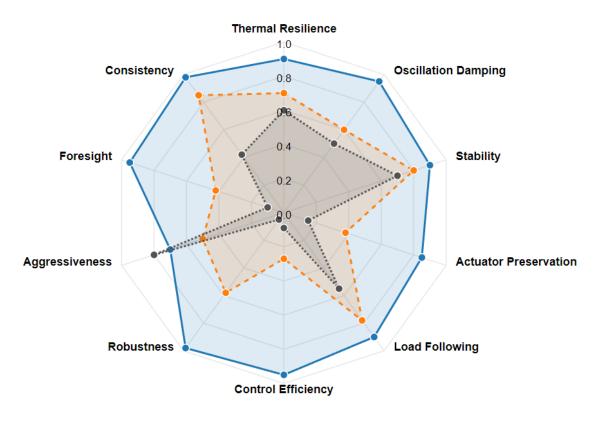
The Scenario-Based Crossover Analysis



The efficiency of the **PID and FLC controllers collapses as the scenario stress increases**, whereas the SAC agent's efficiency remains high. This plot forensically identifies the boundary where an adaptive approach becomes essential.

Integrated Performance Analysis



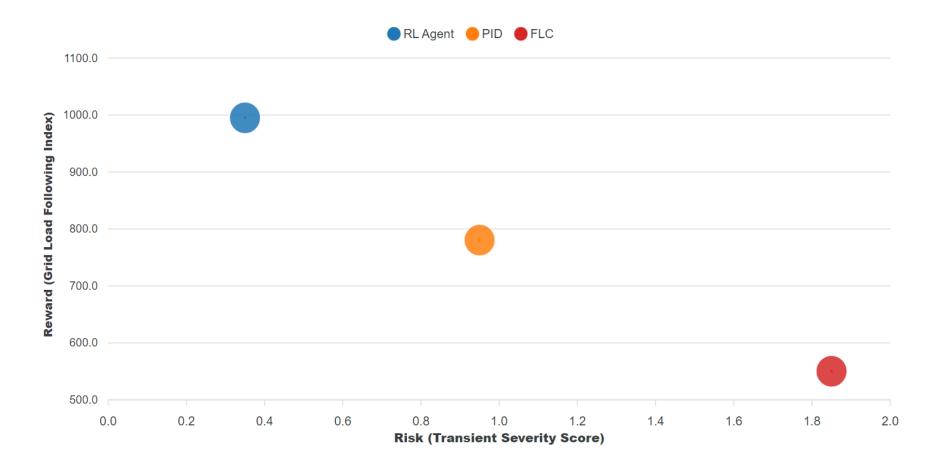


To provide a deeper, qualitative understanding of these aggregate scores, a forensic analysis was conducted on the most demanding scenario that all controllers managed to complete: the *cascading_grid_fault_and_recovery*. The time-series response of the grid frequency which is a critical system state variable.

This the multi-attribute profile gives rise to distinct controller "personalities,"

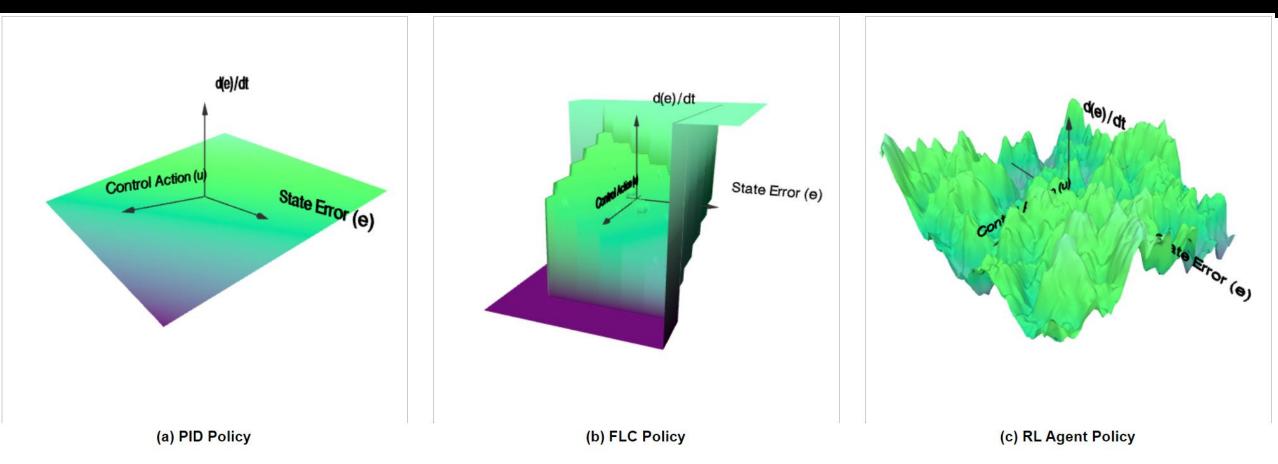
■ RL Agent ● PID ● FLC

Load following Risk during Transient

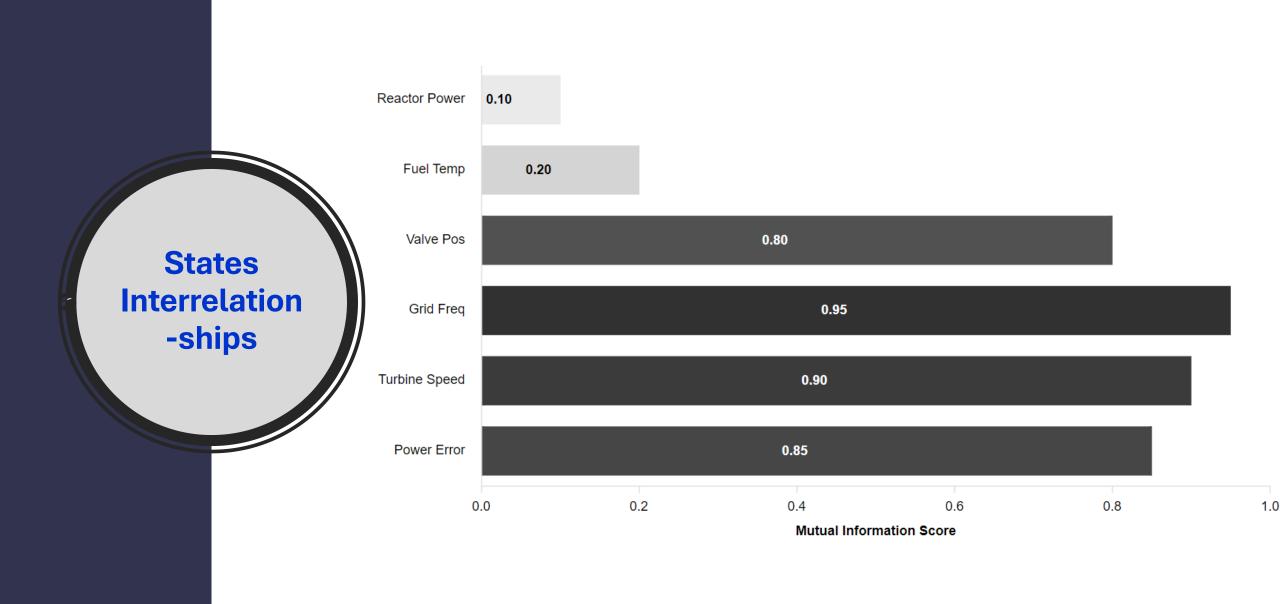


□ This plot illustrates the risk-reward trade-offs. The ideal controller is in the top-left (low risk, high reward). □ The RL Agent clearly operates on the optimal Pareto frontier, achieving the best possible outcome. The classical controllers are in a suboptimal region.

Policy Manifold



- ✓ The foundational reason for the SAC agent's superior performance lies in the distinction between a controller's low-dimensional parameter space and its high-dimensional policy space.
 - ✓ A classical PID controller, even when globally optimized via Differential Evolution, is ultimately defined by *a simple linear plane* in the state-action space,







A New Paradigm for Al Assurance in Critical Systems

The Verdict:
Quantitative
Dominance

25%

Reduction in Control Effort vs. Optimized PID

22%

Superior Composite Robustness Score (p < 0.001)

The Strategic Insight: Architectural Superiority

The Policy Manifold

Classical controllers optimize a **simple policy structure**. The SAC agent discovers a fundamentally **superior**, **high-dimensional policy manifold**.







Scientific & Practical Impact

Innovation

First application of safety-certifiable RL with fuzzy rewards for PWR governor control.

One of Safety

Explicit incorporation of safety constraints aligned with nuclear standards (IAEA, IEEE).

Improved Grid Integration

Enables PWRs to participate more effectively in grids with high renewable penetration.

Impact Pathways for the Nuclear Community:

- Potential framework for Al certification in safety-critical systems.
- Supports **operational flexibility** and modernization efforts in existing and future plants.
- Provides a research benchmark for academic and industrial R&D.

Future Work & Potential Enhancements:











01

Advanced Uncertainty Quantification:

Incorporate deeper Monte Carlo simulations to assess robustness against parameter uncertainties. 02

Formal Safety Verification:

Apply formal methods or Probabilistic Risk Assessment (PRA) techniques for deeper safety analysis (beyond current scope). 03

Hardware-in-the-Loop (HIL) Testing: Bridge the gap between simulation and reality by testing the controller with real hardware components.

04

Expanded Scenarios:

Include more complex grid events or cyberphysical attack scenarios. 05

Explainable AI (XAI):

Develop methods to better understand the RL agent's decisionmaking process.

Thank You! Questions?

ahmedrahman299@gmail.com

