

# Power Grid RL Environment Design

## 1. Environment Components

- **N main generators**: each auto-scales between min & max output.
- **N emergency generators** (bundled): agent can boot or shut them down.
- **1 battery**: can discharge or recharge (charging adds to demand).
- **1 high-priority zone**
- **2 low-priority zones**: each with different scaled demand curves.
- **24-hour dynamic demand cycle**, repeating every episode.
- The agent operates in a **continuous time-step** simulation (e.g., hourly steps) and must keep demand met while minimizing cost and penalties.
- **Demand** is based on a shared Gaussian peak function, scaled per zone.

## 2. Load Zone Demand Functions

Each zone's demand over time  $t \in [0, 24]$  (in hours):

```
def demand_profile(t, scale=1.0):  
    base = 0.5  
    morning_peak = 0.4 * np.exp(-(t - 8)**2 / (2 * 2**2))  
    evening_peak = 0.6 * np.exp(-(t - 19)**2 / (2 * 2**2))  
    return scale * (base + morning_peak + evening_peak)  
  
D_hi = demand_profile(t, scale=1.0)  
D_lo1 = demand_profile(t, scale=0.7)  
D_lo2 = demand_profile(t, scale=0.5)
```

## 3. Generator Behavior

Each **main generator i**:

- Output auto-scales to meet demand **within limits**:

- $G[i]_{\min} \leq \text{output} \leq G[i]_{\max}$
- May fail with probability  $p_{\text{fail}}$  per step:
  - Goes offline and heals after  $T_{\text{down}}$  steps.

Each **emergency generator i** (bundled with main gen):

- Controlled by agent:
  - Can be booted (  $\text{start\_time} = T_{\text{boot}}$  )
  - Can be **shut down manually** to preserve runtime
- When online:
  - Produces fixed power or is controllable
  - Has limited  $\text{runtime\_remaining}$  , decremented when running

## 4. Battery Dynamics

- $\text{SoC} \in [0, \text{max\_capacity}]$
- If agent **discharges**:
  - Draws power, SoC decreases
- If agent **charges**:
  - SoC increases
  - Charging **adds to total demand**, so generator/battery must cover it

Battery control options:

- $a_{\text{batt}} \in \{0: \text{idle}, 1: \text{discharge}, 2: \text{charge}\}$

## 5. Observation Space (Expanded)

Per step, agent sees:

Category	Variables
Time	$t_{\text{norm}} \in [0,1]$
Battery	$\text{SoC}$ , $a_{\text{batt\_mode}}$
Load	$D_{\text{hi}}$ , $D_{\text{lo1}}$ , $D_{\text{lo2}}$
Generator[i]	$G[i]_{\text{online}} \in \{0,1\}$ , $G[i]_{\text{fail\_timer}}$ ,

	$G[i]_{\min}$ , $G[i]_{\max}$
Emergency[i]	$E[i]_{\text{online}}$ , $E[i]_{\text{start\_timer}}$ , $E[i]_{\text{runtime\_left}}$

Observation dimension:

$$1 \text{ (time)} + 3 \text{ (loads)} + 1 \text{ (battery)} + N \times (3 + 3) = 4 + 6N$$

## 6. Action Space

For each step:

- $a_{\text{batt}} \in \{0,1,2\}$  — idle, discharge, charge
- $a_{\text{em}}[i] \in \{0,1,2\}$  — do nothing, boot emergency, shut down
- $a_{\text{shed\_lo1}} \in \{0,1\}$  — shed zone 1
- $a_{\text{shed\_lo2}} \in \{0,1\}$  — shed zone 2

Total dimension:

$$3 \text{ (battery)} + N \times 3 \text{ (em control)} + 2 \text{ (load sheds)}$$

## 7. Power Allocation Logic

1. **Total demand** =

$$D_{\text{hi}} + D_{\text{lo1}} * (\text{not shed}) + D_{\text{lo2}} * (\text{not shed}) + \text{battery\_charge}(\text{if charging})$$

2. **Available power** =

$$\sum \text{active main outputs} + \sum \text{emergency outputs} + \text{battery\_discharge}(\text{if discharging})$$

3. **Balance:**

- Try to meet demand in this order:
  1.  $D_{\text{hi}}$
  2.  $D_{\text{lo1}}$
  3.  $D_{\text{lo2}}$
  4. battery charge

## 8. Reward Function

At each step:

$$\text{reward} = ( \\ + w_{\text{hi}} * 1[\text{hi\_zone powered}]$$

```
+ w_lo1 * 1[lo1 powered]
+ w_lo2 * 1[lo2 powered]
- c_batt_discharge * energy_drawn
- c_batt_charge * energy_charged
-  $\sum c_{em\_boot} * 1[emergency\ i\ booted]$ 
-  $\sum c_{em\_run} * 1[emergency\ i\ running]$ 
-  $\sum c_{em\_idle} * 1[emergency\ i\ idle\ but\ online]$ 
-  $\sum c_{fail} * 1[main\ gen\ i\ failed]$ 
- c_shed * 1[any load zone shed]
)
```

- Emergency generators have **run cost** and **boot cost**
- Penalty for keeping unused emergency generators online
- Charging battery has a small cost since it adds to demand

## ✅ Ready to Proceed?

Would you like:

1. A **code scaffold for this environment** (e.g., `PowerGridEnv(gym.Env)`)?
2. A **visual flowchart** of how state → action → reward flows?
3. Or begin with **Q-table / PPO agent plan** for this?

Let's move to implementation.