

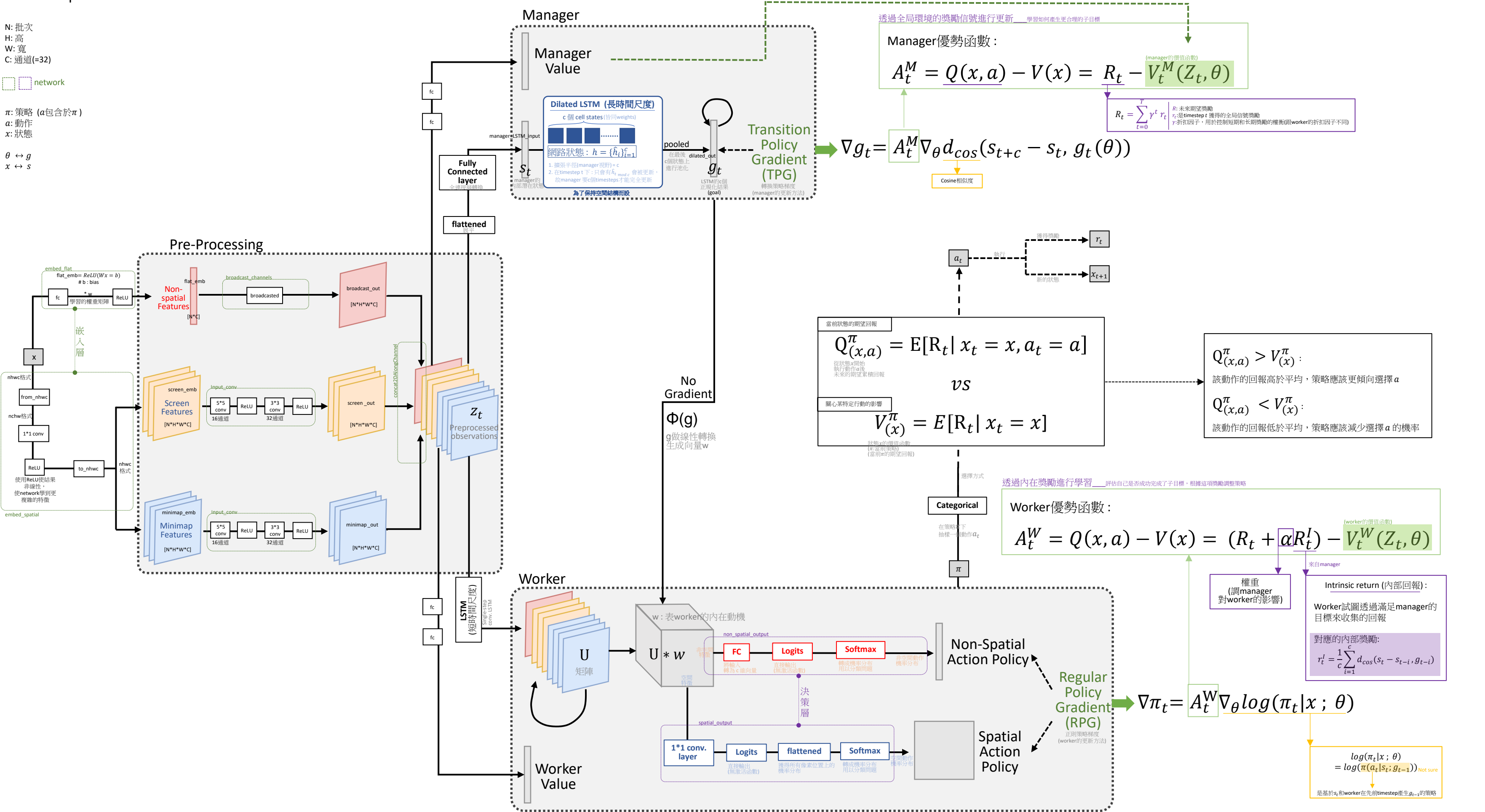
More Complex

N: 批次
H: 高
W: 寬
C: 通道(=32)

network

π : 策略 (a 包含於 π)
 a : 動作
 x : 狀態

$\theta \leftrightarrow g$
 $x \leftrightarrow s$



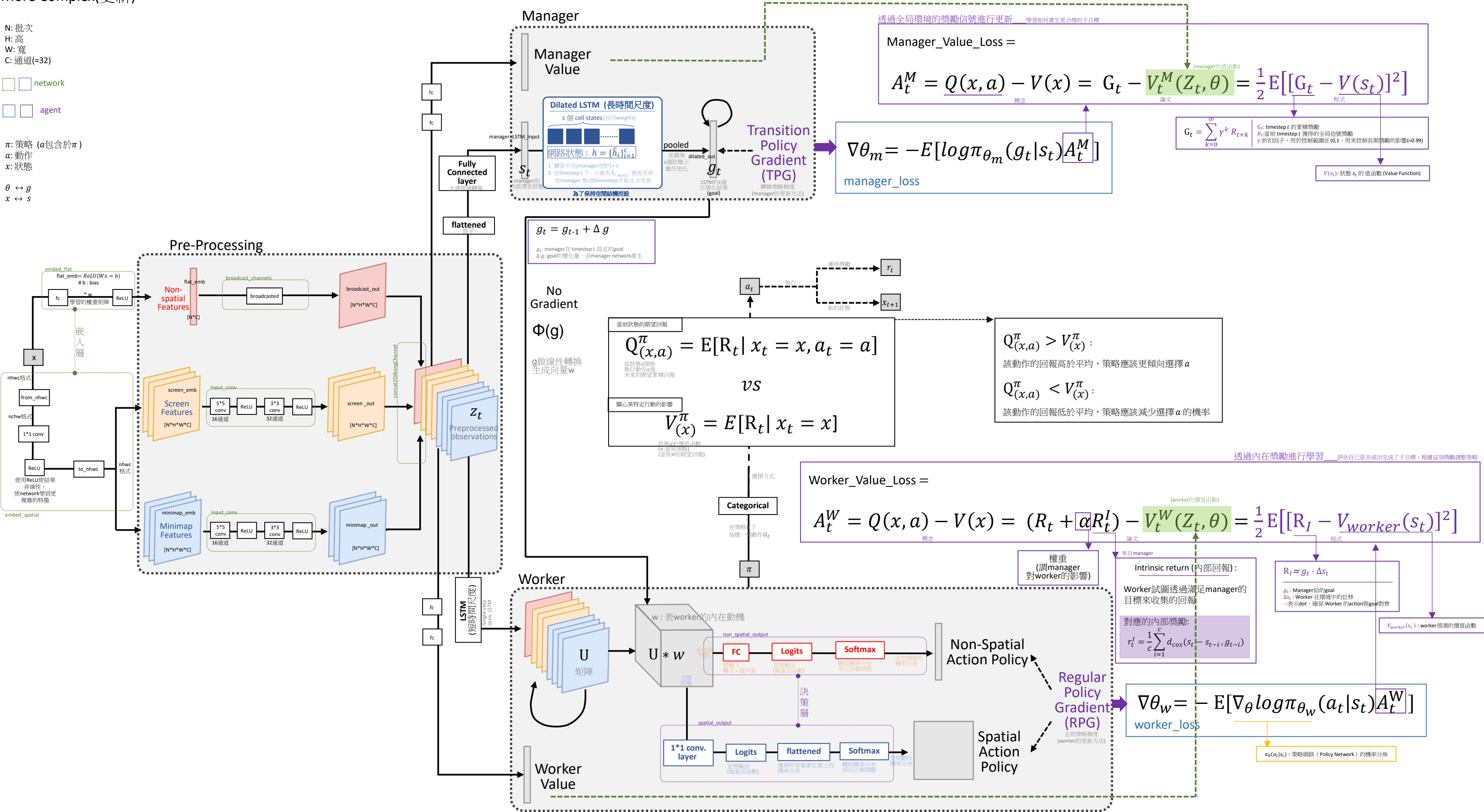
More Complex(更新)

N: 批次
H: 高
W: 寬
C: 通道(=32)

network

agent

π : 策略 (a 包含於 π)
 a : 動作
 x : 狀態
 $\theta \leftrightarrow g$
 $x \leftrightarrow s$



```
python "C:\Users\yuh\u\pysc2-rl-agents\run.py"
experiment_03
--agent feudal
--map MoveToBeacon
--envs 4
--res 32
--steps_per_batch 32
--iters 100000
--lr 0.0005
--entropy_weight 0.01
--save_iters 2000
--summary_iters 50
--vis
--value_loss_weight 0.7
--discount 0.99
```

學習率 α 使用指數衰減：

$$\alpha_t = \alpha_0 \cdot 0.94^{\frac{t}{10000}}$$

- 初始學習率 $\alpha_0 = 0.0005$
- 每 10,000 步，學習率乘以 0.94

$$L_{\text{entropy, manager}} = -\mathbb{E} \left[\sum_g \pi_{\theta_m}(g|s) \log \pi_{\theta_m}(g|s) \right]$$

$$L_{\text{entropy, worker}} = -\mathbb{E} \left[\sum_a \pi_{\theta_w}(a|s) \log \pi_{\theta_w}(a|s) \right]$$

manager_entropy + worker_entropy

$$(\text{manager_loss} + \text{worker_loss}) + \lambda_1 \cdot \text{Value_loss_weight} * (\text{manager_value_loss} + \text{worker_value_loss}) + \lambda_2 \cdot \text{entropy_weight} * (\text{entropy})$$

策略損失 (Policy Loss)

來自 Actor-Critic 方法
目標是最大化累積獎勵

值函數損失 (Value Loss)

用於減少 Agent 預測的價值
和 實際累積獎勵 之間的誤差

熵損失 (Entropy Loss)

熵損失 $H(\pi)$ 用來增加策略的隨機性
以鼓勵探索