



$$f_{\theta}(\tau_f(s_t))$$

$$\tau_g(s_{\text{detox}})$$



$$g_{\text{toxic}}^*(f_{\theta}(\tau_f(s_t))) \approx g_{\text{toxic}}(\tau_g(s_{\text{detox}}))$$