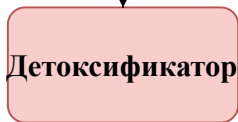
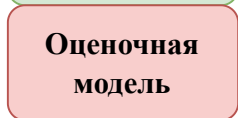
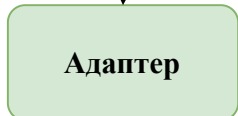


$$\tau_f(s_t)$$

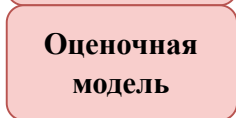
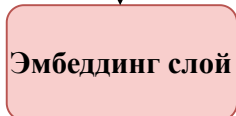


$$f_{\theta}(\tau_f(s_t))$$



$$P_{\text{toxic}}^*$$

$$\tau_g(s_d)$$



$$P_{\text{toxic}}$$

 g_{toxic}

 g_{toxic}^*


$$D_{\text{KL}}(P_{\text{toxic}} \parallel P_{\text{toxic}}^*)$$