



Figure 5: t-SNE visualization of harmful embedding drift under different harmful ratios p . Each point represents the embedding of each alignment data. As shown, for SFT, when harmful ratio p is high, we observe that the hidden embedding of the alignment data drifts significantly from that of the model before finetuning. For Vaccine, we observe a mitigated embedding drift, which explains why Vaccine is still able to maintain alignment knowledge.

Table 10: Performance evaluation of Accelerated Vaccine. As shown, Accelerated Vaccine with proper perturbation periodicity can maintain defense performance, while significantly reducing the training step time.

Methods	Harmful score	Finetune accuracy	Training time
SFT	59.20	94.40	0.12902s (1x)
Vaccine	50.40 (+0)	92.40	0.24852s (1.93x)
Accelerated Vaccine($\tau = 100$)	51.00 (+0.60)	95.20	0.13140s (1.02x)
Accelerated Vaccine($\tau = 1000$)	52.00 (+1.60)	94.20	0.12956s (1.004x)
Accelerated Vaccine($\tau = 10000$)	53.20 (+2.80)	94.80	0.12934s (1.002x)
Accelerated Vaccine($\tau = 20000$)	58.80 (+8.40)	94.40	0.12902s (1x)

Algorithm 2 Accelerated Vaccine (Perturb every τ steps)

input Perturbation intensity ρ ; Local step T ; Layer number L ; **Perturbation Periodicity** τ ;
output The aligned model w_{t+1} ready for finetuning.
for step $t \in T$ **do**
 Sample a batch of data (x_t, y_t)
 Backward $\nabla \mathcal{L}_{w_t}(e_{1,t}, \dots, e_{L,t})$ with (x_t, y_t)
 if $t \bmod \tau == 0$ **then**
 for each layer $l \in L$ **do**
 $\epsilon_{l,t} = \rho \frac{\nabla_{e_{l,t}} \mathcal{L}_{w_t}(e_{l,t})}{\|\nabla \mathcal{L}_{w_t}(e_{1,t}, \dots, e_{L,t})\|}$
 Register forward hook: $\tilde{f}_{w_l, \epsilon_{l,t}}(e_{l,t}) = f_{w_l}(e_{l,t}) + \epsilon_{l,t}$
 end for
 Backward $\tilde{g}_t = \nabla \mathcal{L}((\tilde{f}_{w_{L,t}, \epsilon_{L,t}} \circ \dots \circ \tilde{f}_{w_{1,t}, \epsilon_{1,t}}) \circ \mathcal{T}(x_t, y_t))$
 else
 $\tilde{g}_t = \nabla \mathcal{L}_{w_t}(e_{1,t}, \dots, e_{L,t})$ with (x_t, y_t)
 end if
 $w_{t+1} = \text{Optimizer_Step}(w_t, \tilde{g}_t)$
end for
