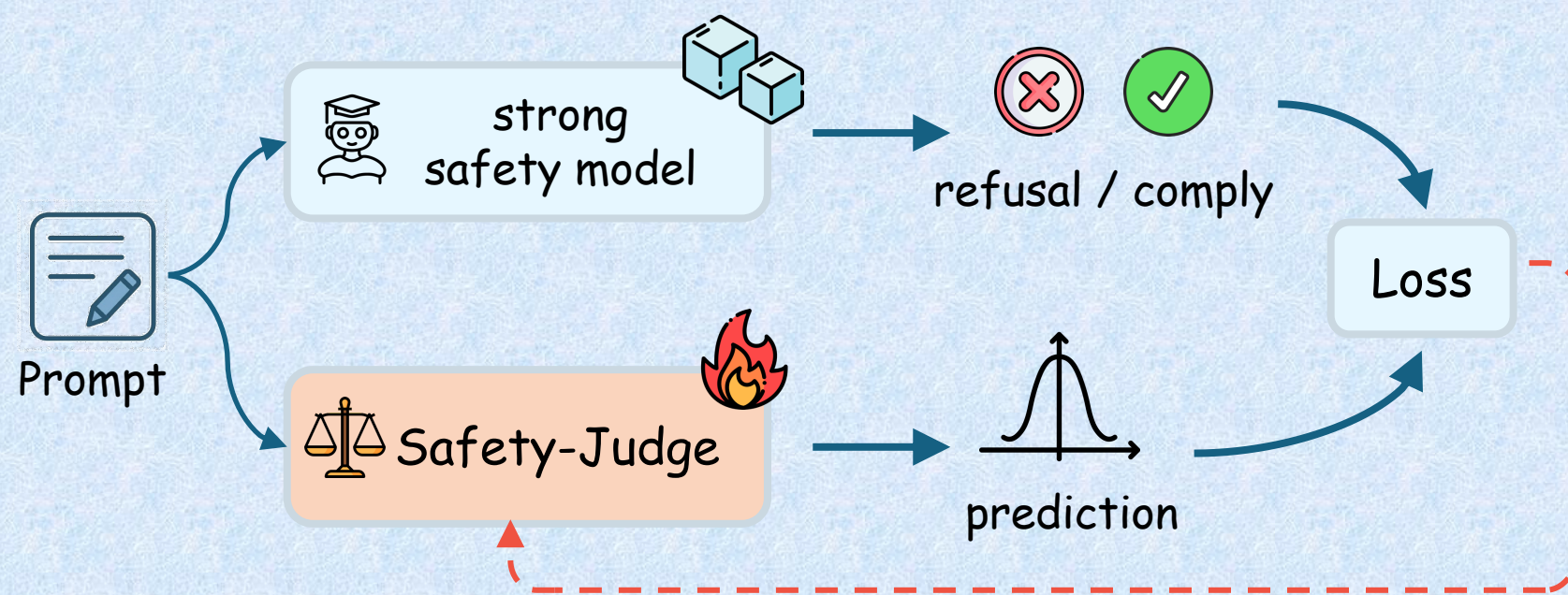
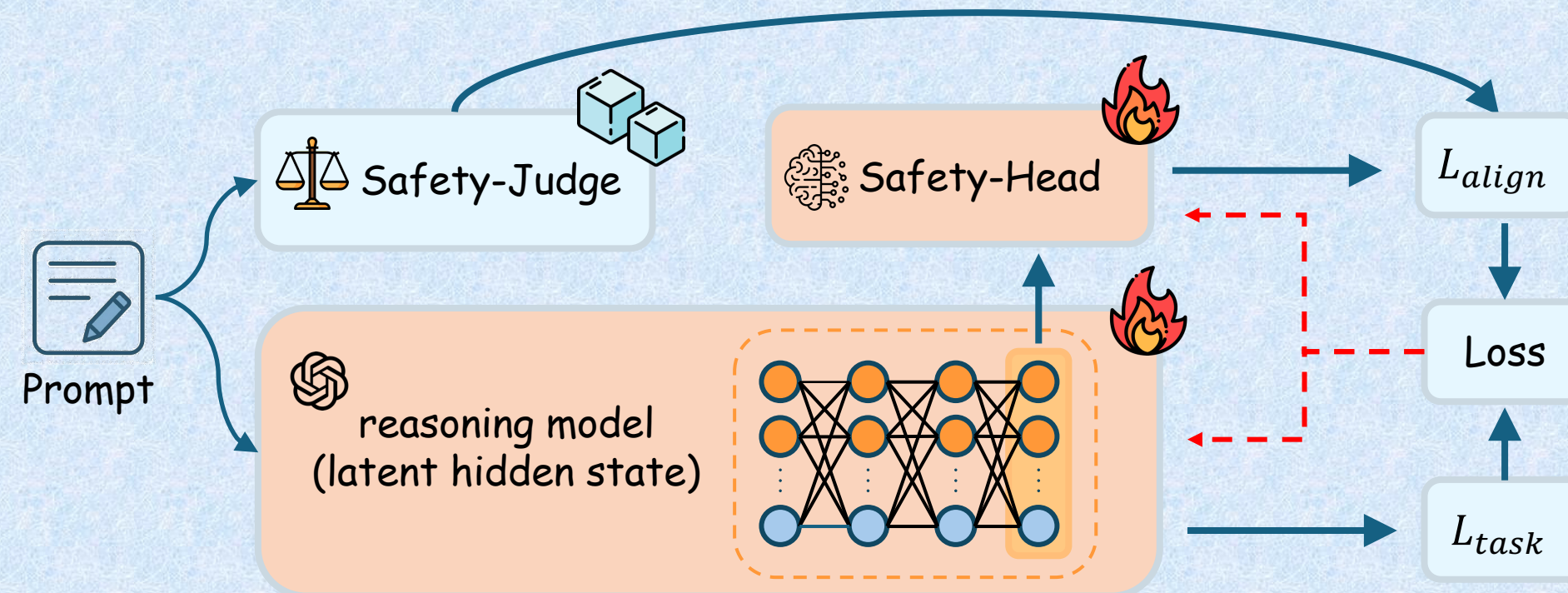


Training phase

safety policy distillation



latent space safety alignment



Inference phase

