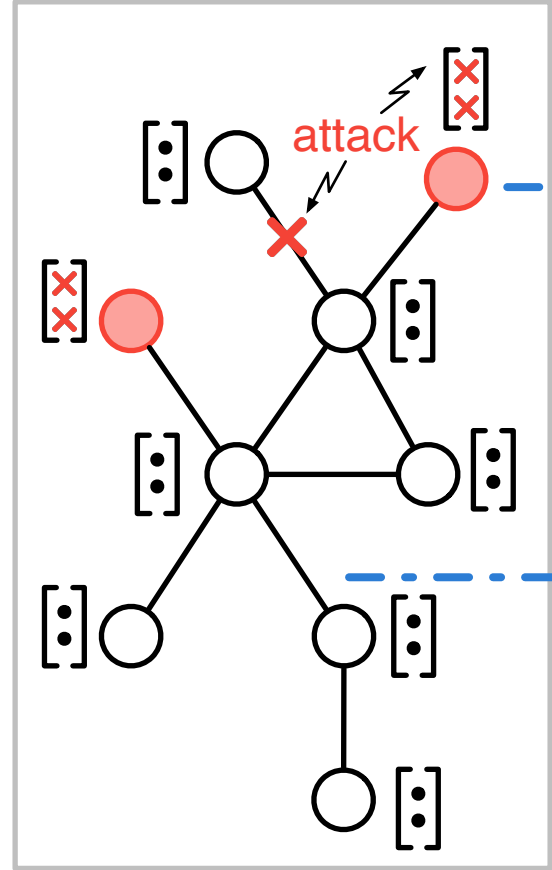




Graph modification attack

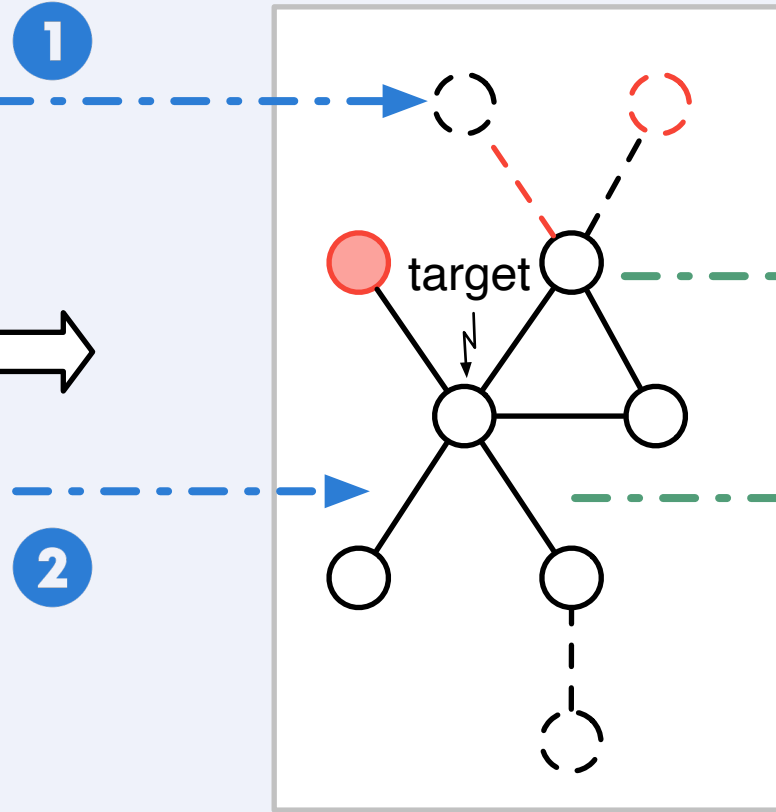
○ benign ○ attacked



adversarial graph
 $\tilde{G} = (\tilde{X}, \tilde{A})$

1 Subgraph Extractor

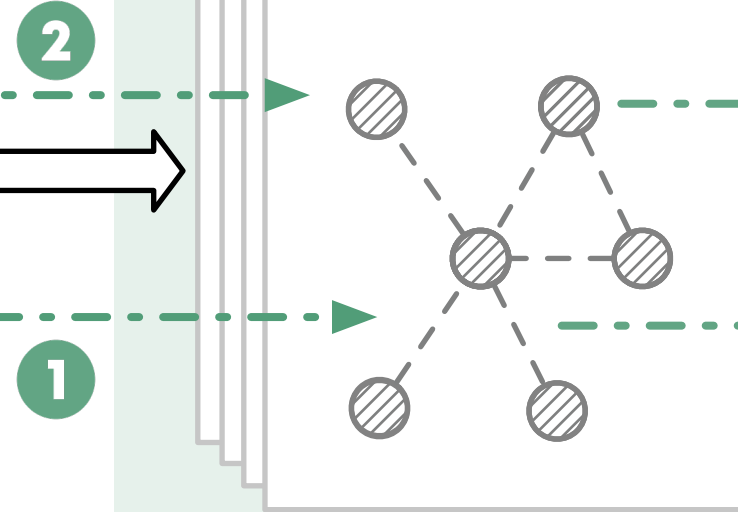
- 1 $X' = \{u \mid d(u, v) \leq L\}$
- 2 $A' = \{(u, w) \mid u, w \in X'\}$



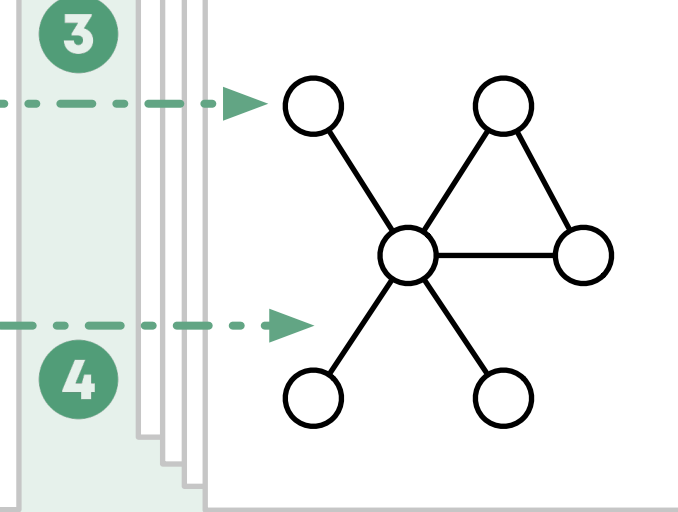
subgraph $G' = (X', A')$ of target v

2 Denoising Graph Joint Diffuser

- 1 $A_t \sim A' \oplus \epsilon_{t_a}^A$
- 2 $X_t \sim X' \oplus \epsilon_{t_x}^X$
- 3 $\bar{X} \sim X_t \oplus \epsilon_{\theta}(X_t, t_x)$
- 4 $\bar{A} \sim A_t \oplus \epsilon_{\theta}(A_t, t_a \mid \bar{X})$



diffused graphs
 $G_t = (X_t, A_t)$



denoised graphs
 $\bar{G} = (\bar{X}, \bar{A})$

3 Certificate Generator



robustness certificate

$$f_{\theta}(v \mid \tilde{G}) = y^*, \quad ||\tilde{G} - G|| \in \Delta$$



GNN $f_{\theta}(v \mid G)$

n predictions

