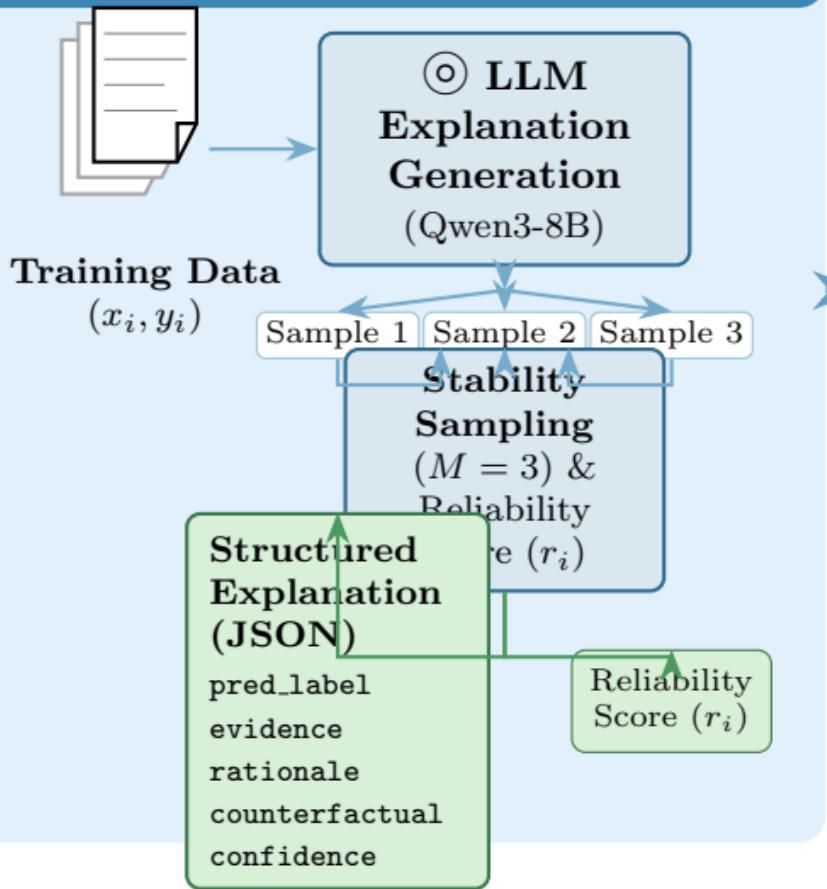
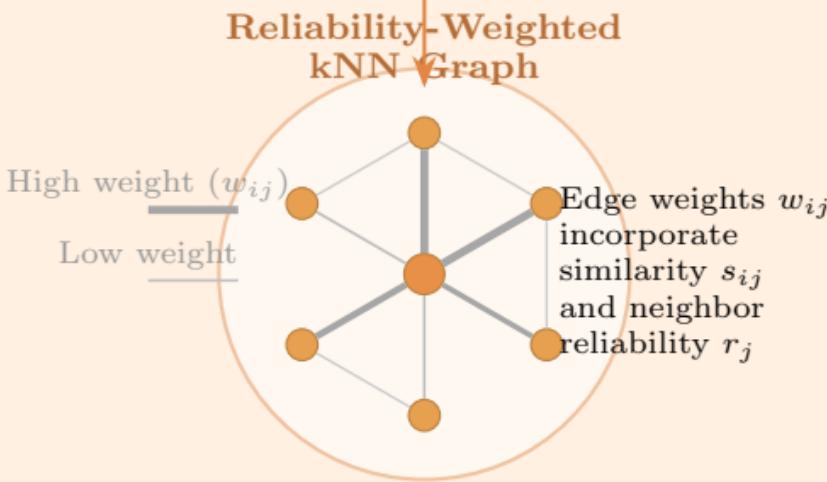
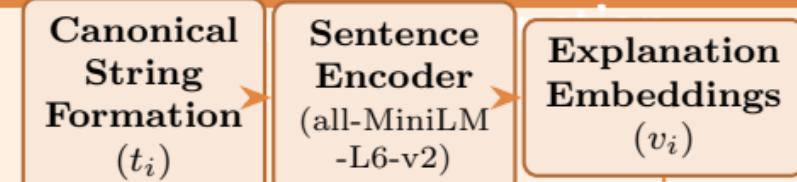


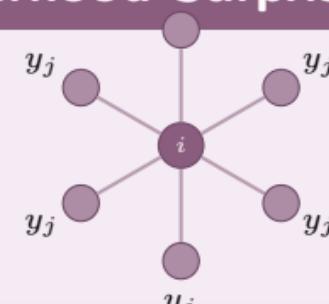
1. Explanation Generation



2. Explanation Embedding



3. Neighborhood Surprise Detection



Weighted neighbor label posterior

$$p_i(c) = \sum_j w_{ij} \cdot \mathbf{1}[y_j = c]$$

Neighborhood Surprise (S_n)
 $= -\log p_i(y_i)$

Suspiciousness Ranking
Flagged Instances

