**Token Scores = Base Loss − Reference Loss**

**Fixed-Model Cleaning**

**Token Scores for $\tilde{D}_t$ = Base Loss − t-th Reference Loss**

**Self-Evolving Cleaning**