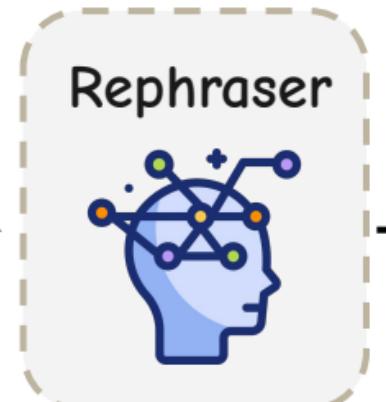


Rephraser Training

Organic Data



Recycled Data

RL Training

Faithfulness Rewards
Semantics, Structure, Length

Quality Reward
DataMan

Web Recycling

Organic Data Pool



Quality Filtering



Final Pretraining Corpus



Recycled Data Pool

