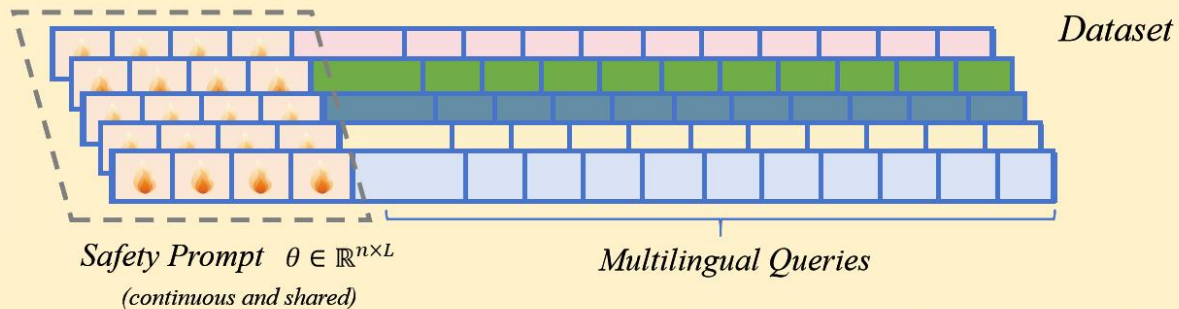
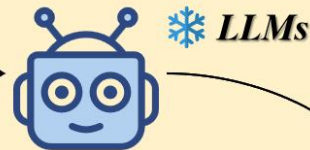
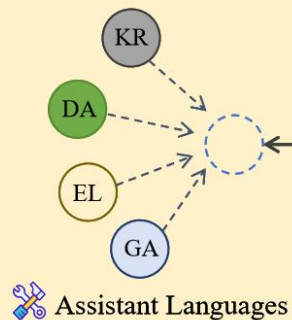


### Initial Safety Prompt (EN)

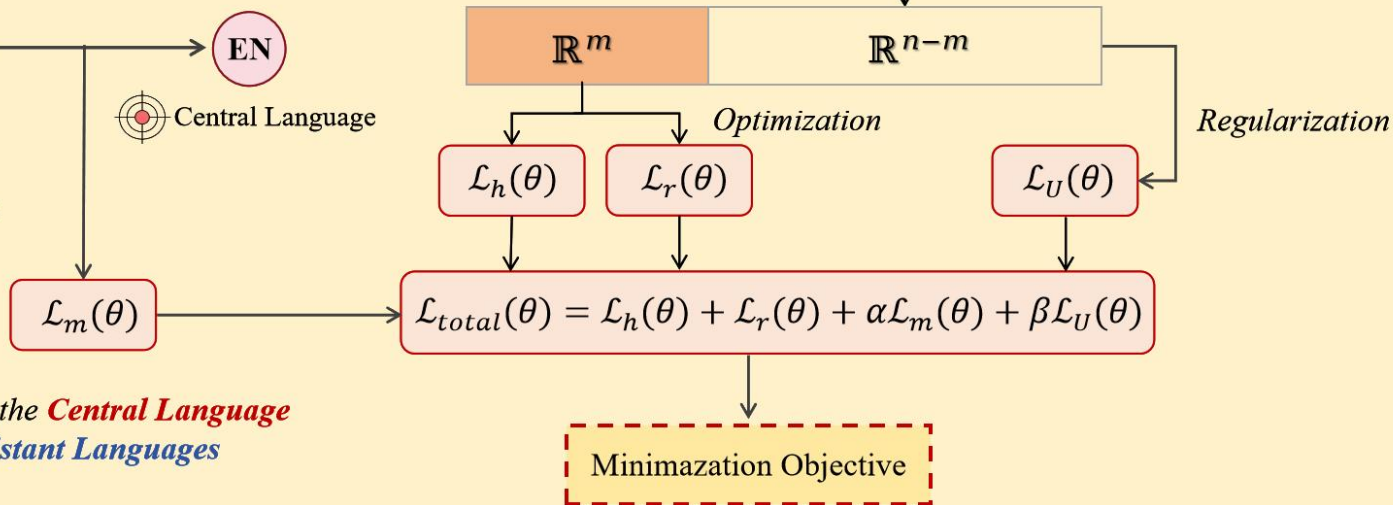
You are a **helpful respectful and honest** assistant . Always answer as helpfully as possible, while being safe.



### Multilingual Collaborative Collection



- $\mathcal{L}_h(\theta)$  : Harmful Query Distinction
- $\mathcal{L}_r(\theta)$  : Refusal Probabilities Optimization
- $\mathcal{L}_U(\theta)$  : N-M Dimensional Regularization
- $\mathcal{L}_m(\theta)$  : Multilingual Safety Alignment



Set a language as the **Central Language** and others as **Assistant Languages**