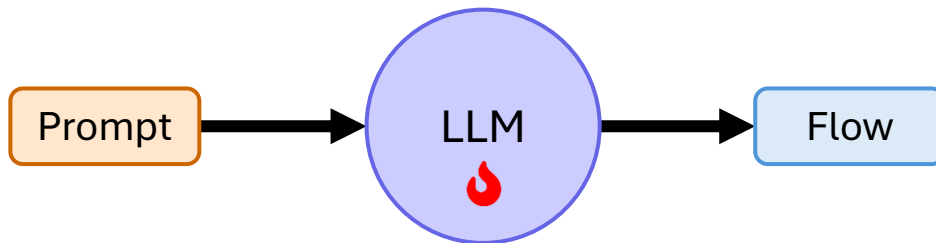


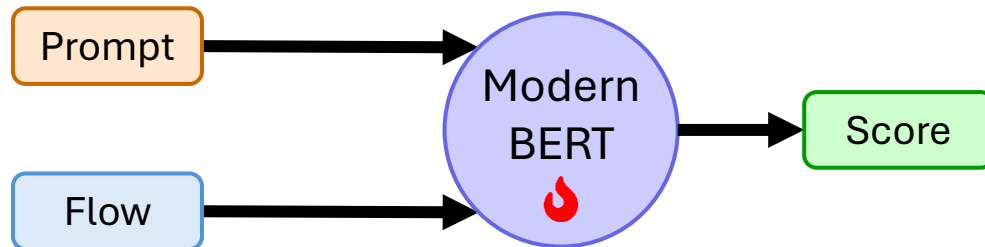
Stage 1: SFT

Prompt $x500k$
Flow pairs
Data



Stage 2.1: Surrogate Reward Training

Prompt **Score**
Flow $x100k$ triplets
Data



Stage 2.2: Policy-Optimization

Prompt $x5000$
Data

