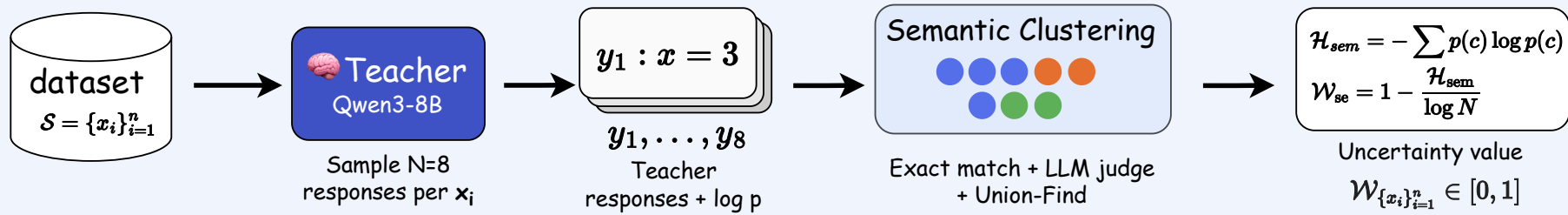


## Phase 1: Offline Uncertainty Estimation No training cost



precomputed  $\mathcal{W}_{se}$  stored with data

## Phase 2: Online Uncertainty-Calibrated Training On-policy loop

