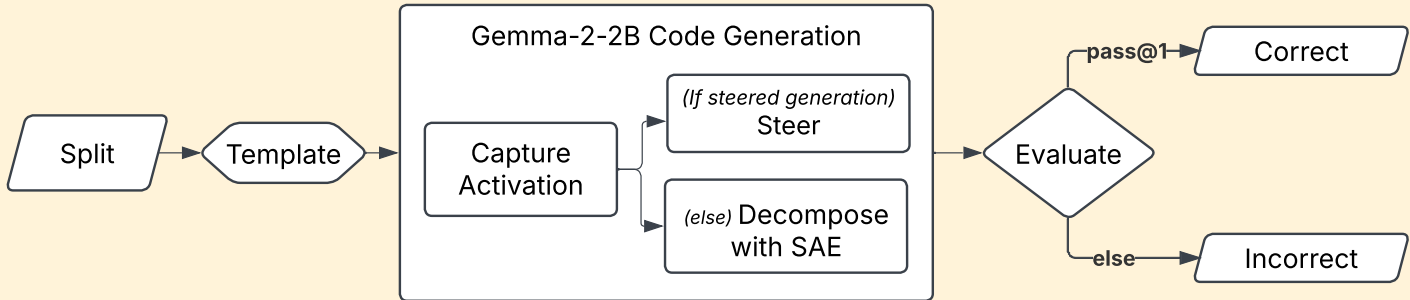
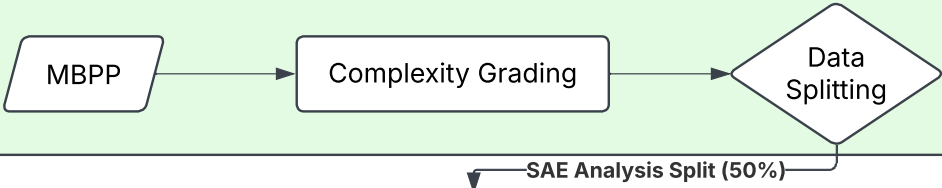


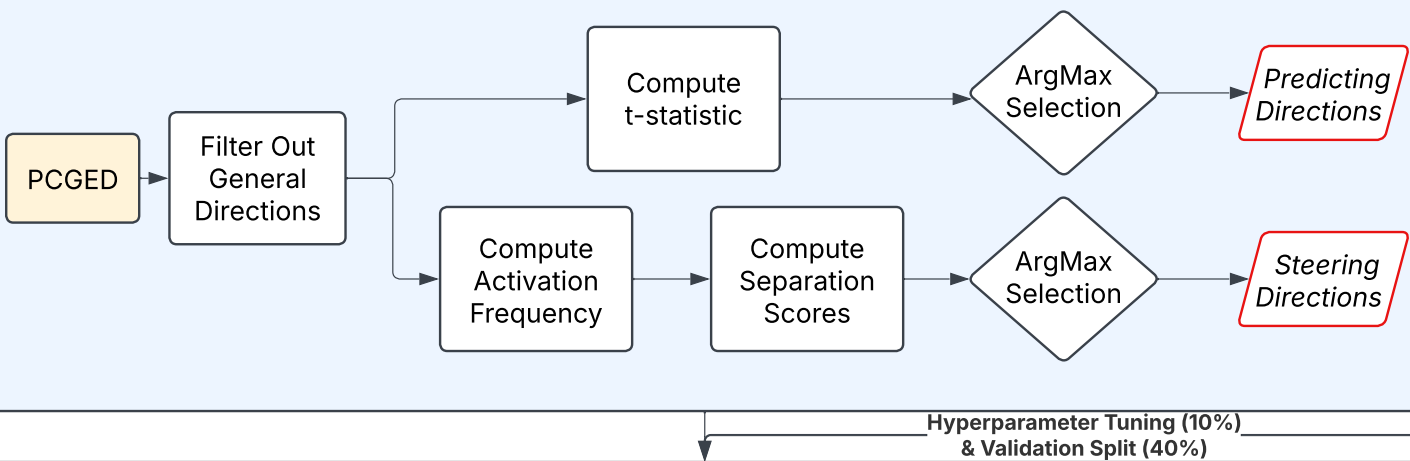
Prompt-Capture-Decompose-Generate-Evaluate (PCDGE)



1. Dataset Preparation

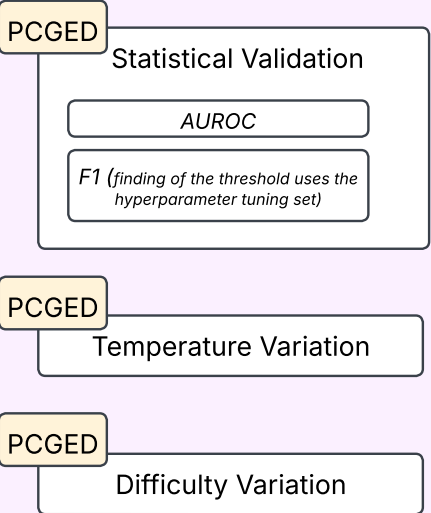


2. Direction Selection

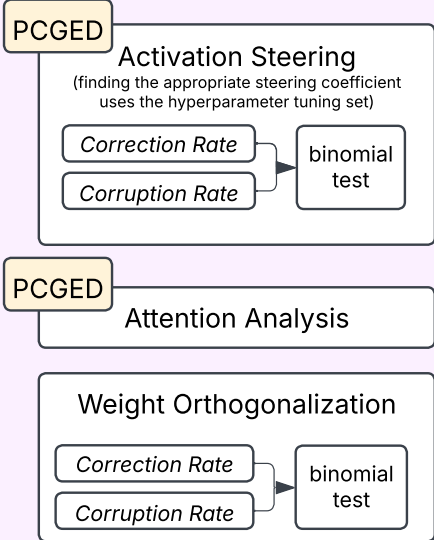


3. Mechanistic Analysis

Predicting Direction



Steering Direction



Logit Lens Analysis

Persistence Testing