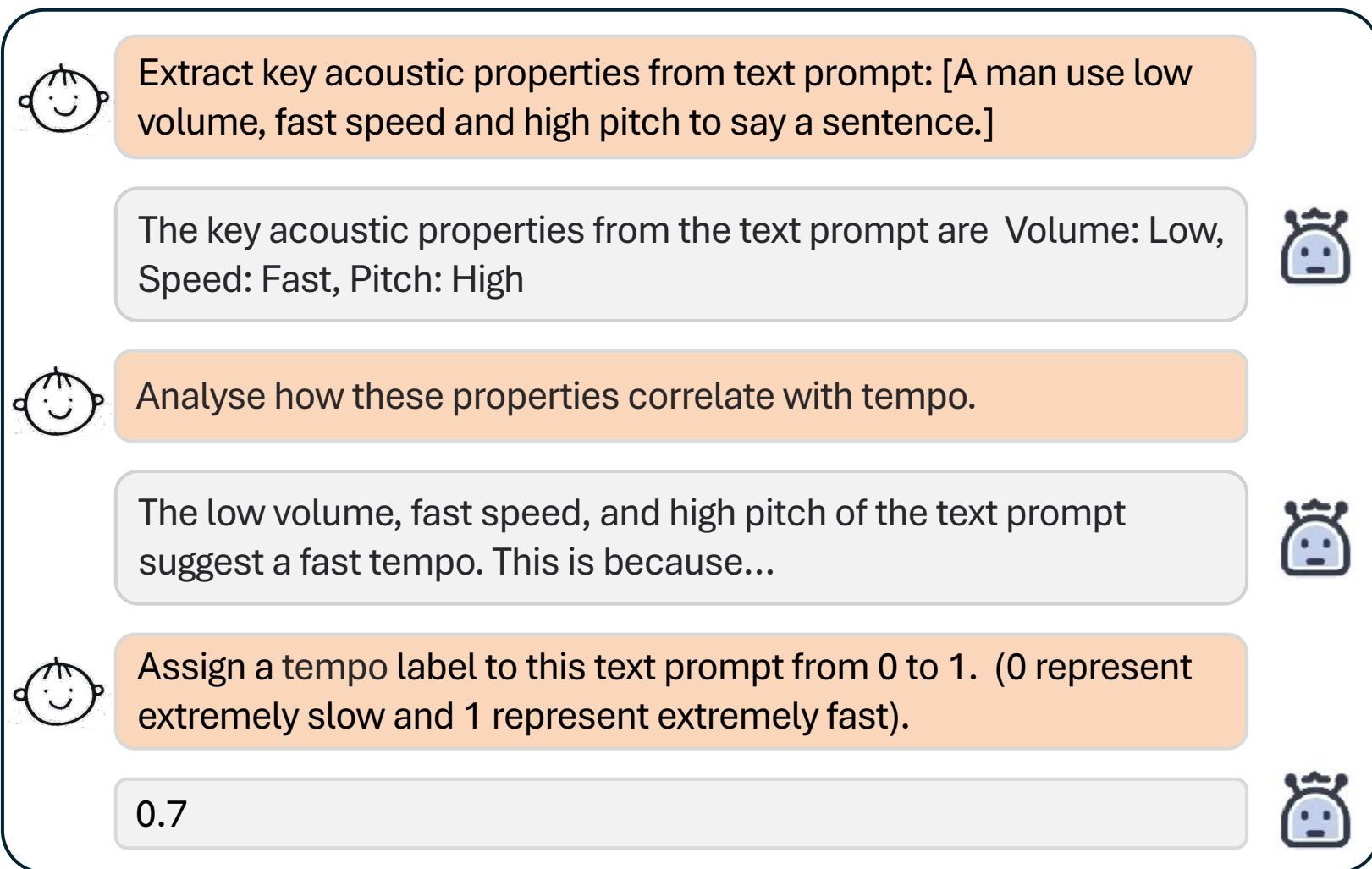
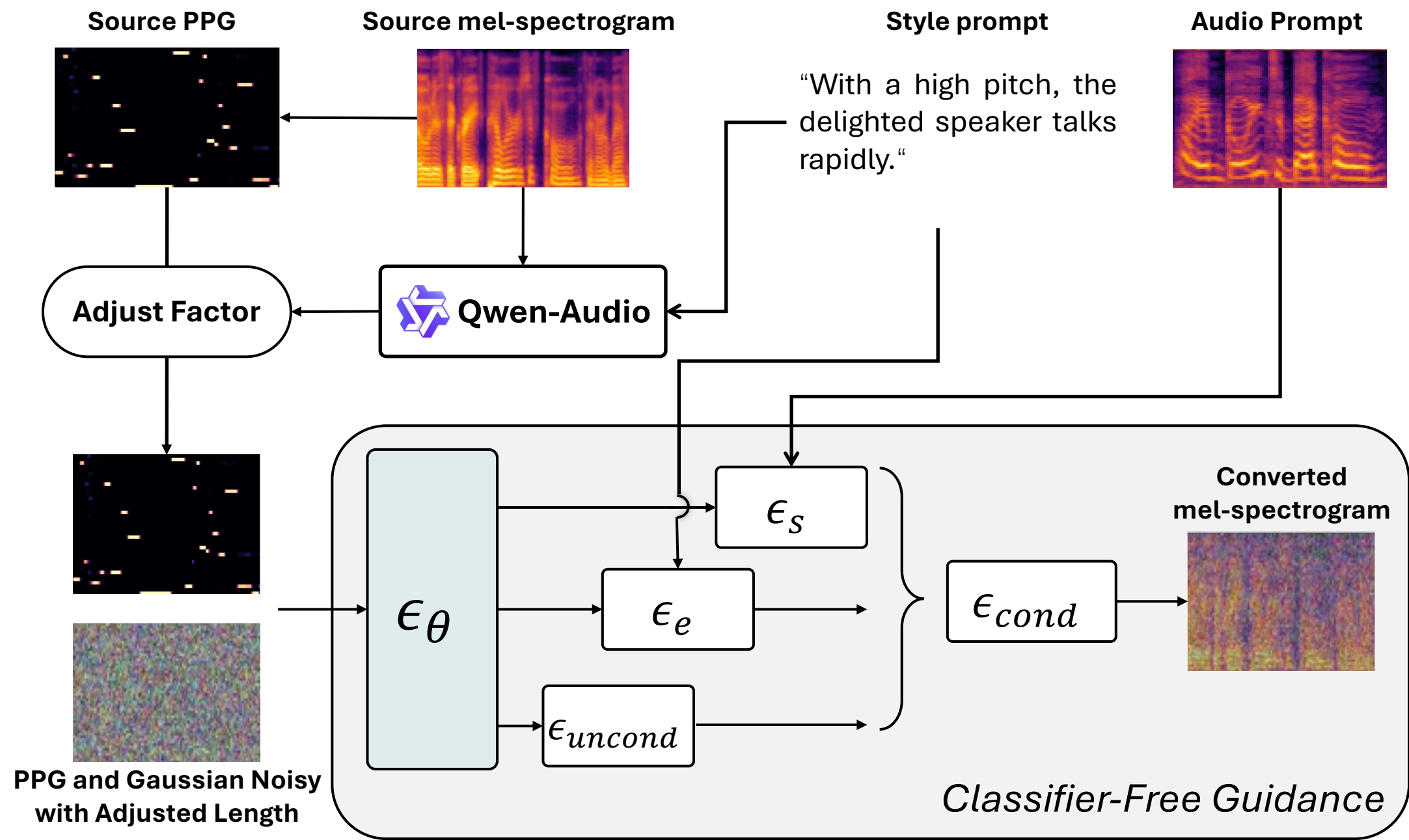


(a) Tempo estimation from source speech



(b) Tempo estimation from text prompt



(c) Inference stage of the proposed UniVST