Figure 1: Sampling algorithms choose the next token at each time step $s_i$ by sampling from the conditional distribution $p_\theta(\cdot \mid \boldsymbol{y}_{:i-1})$ and appending it to the context.