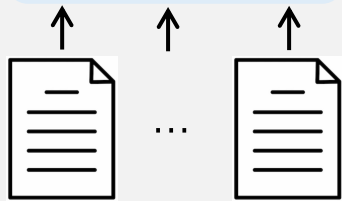


Score & human
↑ feedback

LLMs



y_1

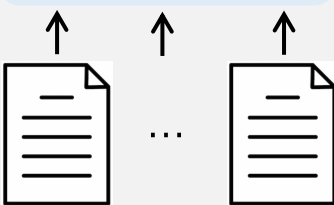
y_m

(a) benchmark paradigm

Judge



LLMs



y_1

y_m

(b) LLMs as judges

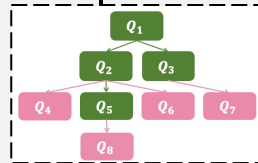
Judge



LLMs



Examiner



(c) TreeEval