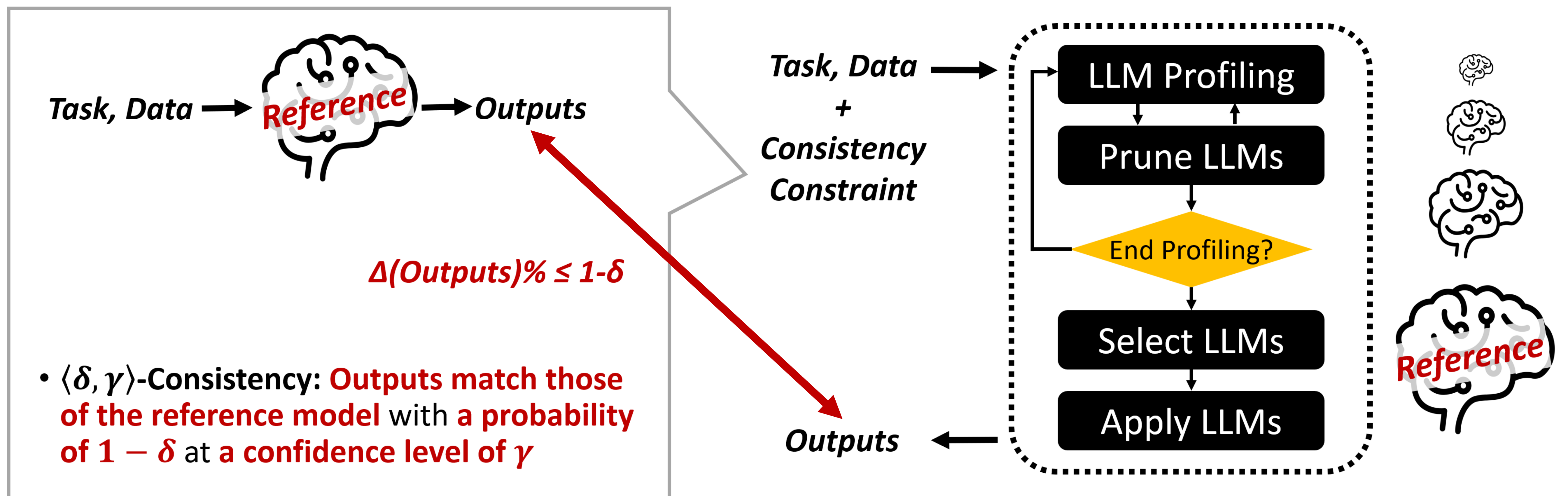


SpareLLM

Automatically Selecting Task-Specific Minimum-Cost Large Language Models under Equivalence Constraint

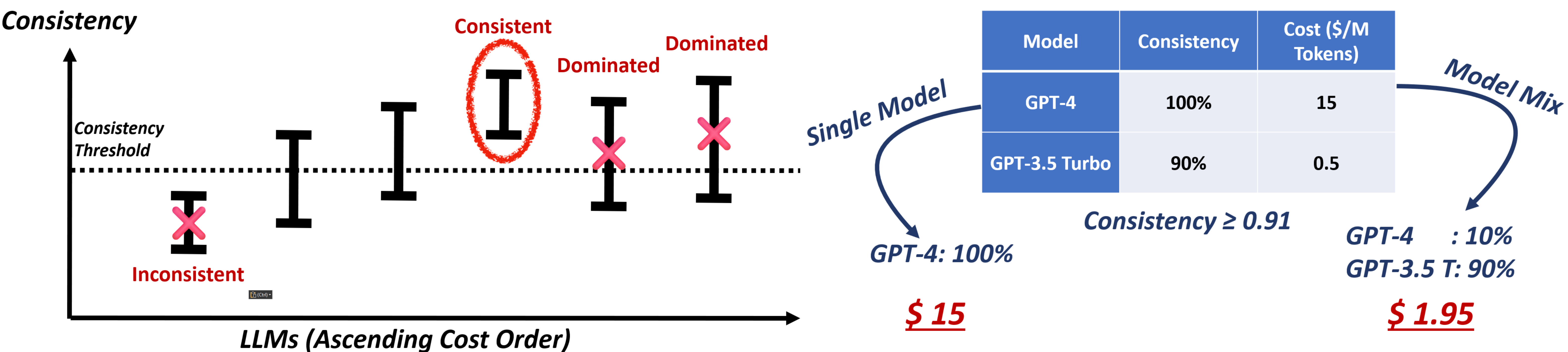
Saehan Jo and Immanuel Trummer
Cornell University

SpareLLM Architecture



Profiling Phase

Application Phase



Experiments

