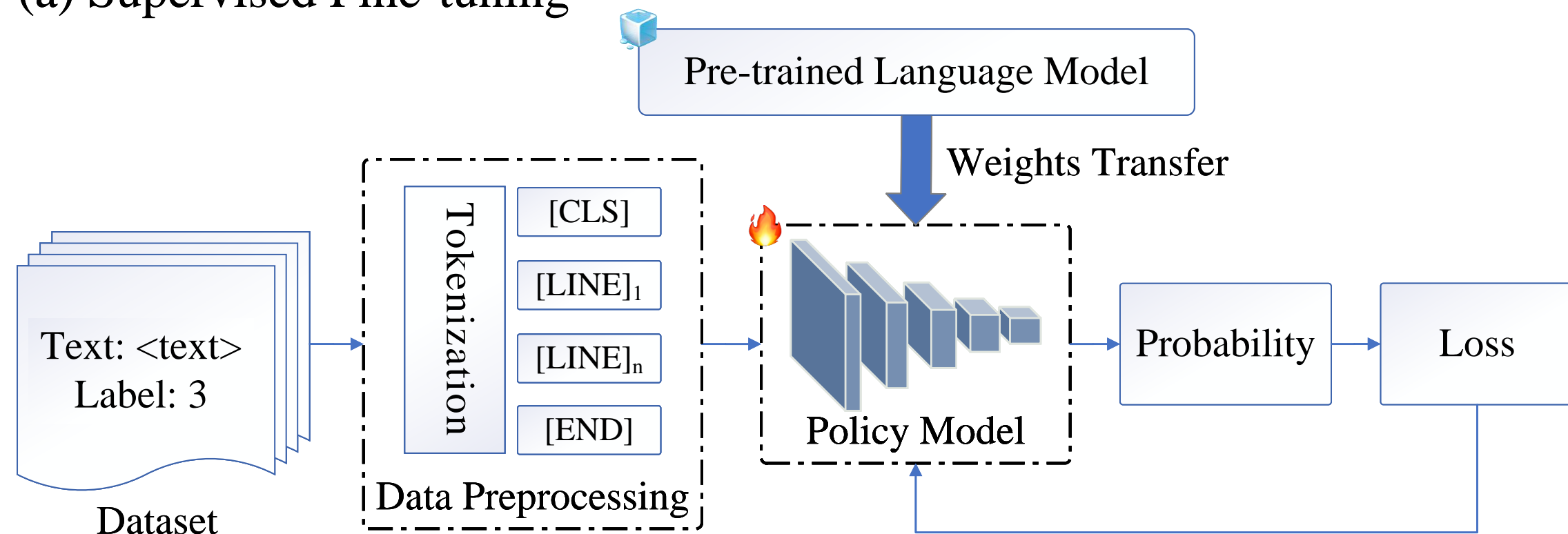
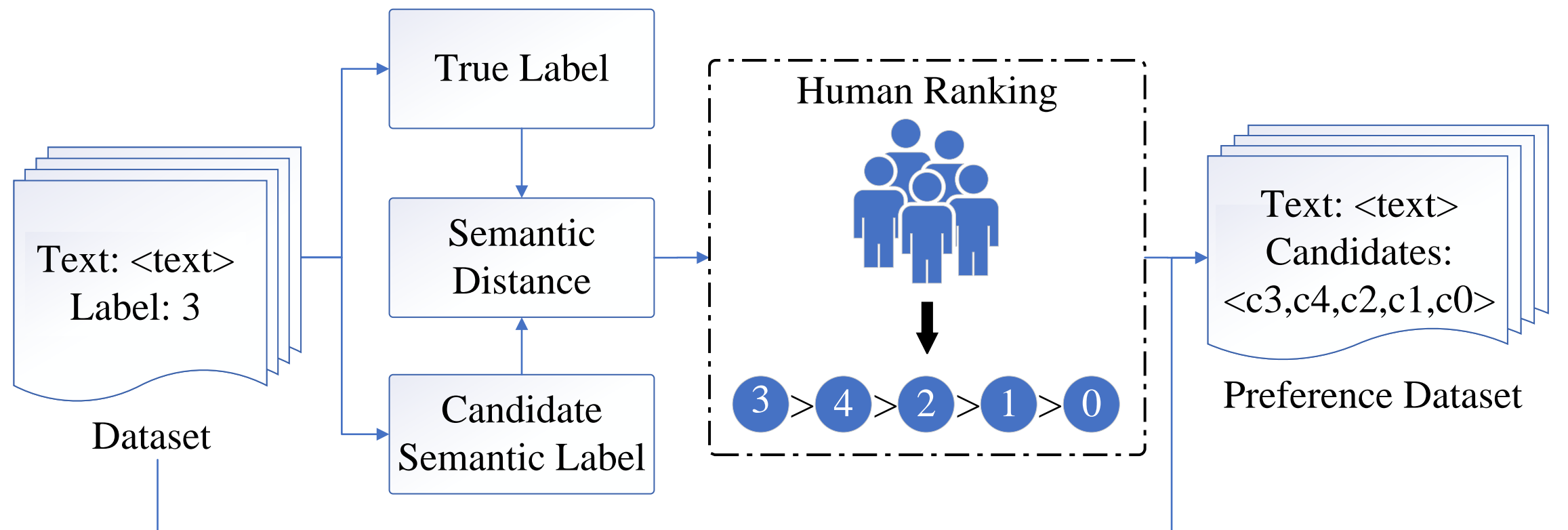


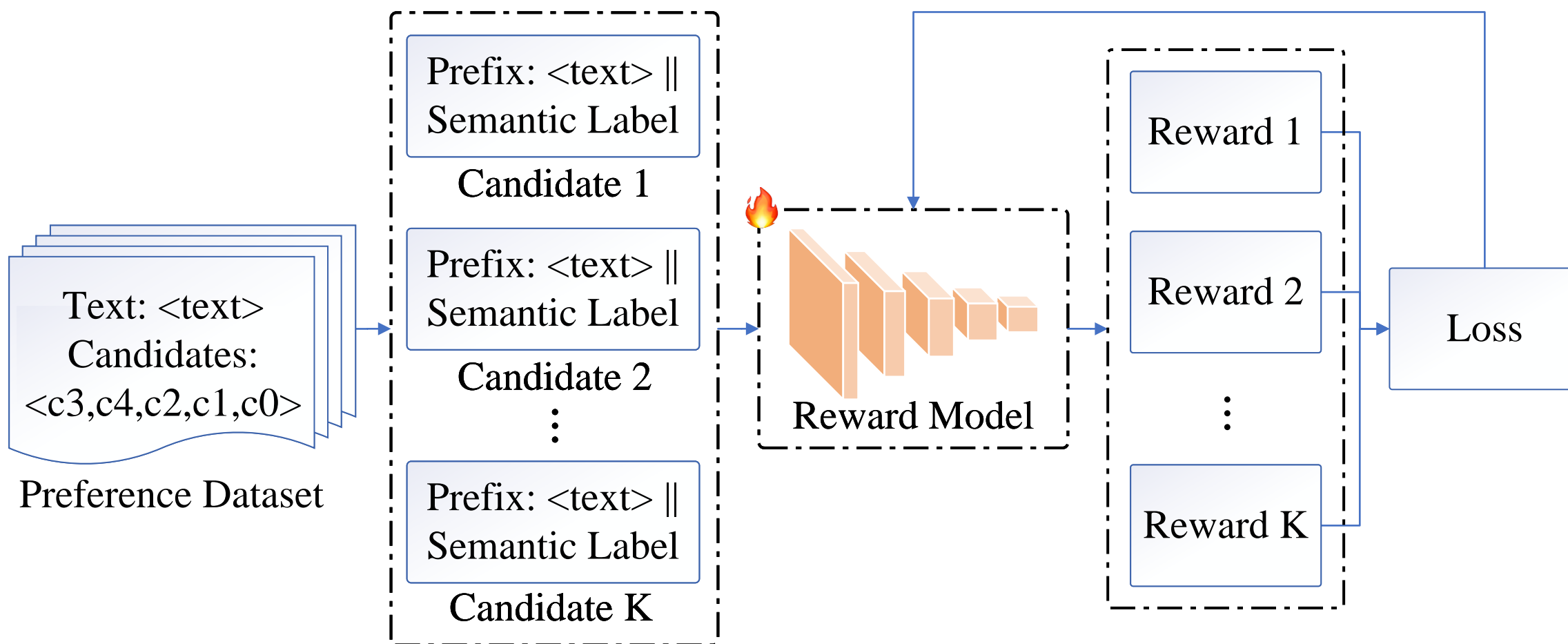
(a) Supervised Fine-tuning



(b) Preference Data Construction



(c) Reward Model Training



(d) Reinforcement Learning Optimization

