

# Affordable On-line Dialogue Policy Learning

Cheng Chang\*, Runzhe Yang\*, Lu Chen, Xiang Zhou and Kai Yu

Key Lab. of Shanghai Education Commission for Intelligent Interaction and Cognitive Eng.

SpeechLab, Department of Computer Science and Engineering

Brain Science and Technology Research Center

Shanghai Jiao Tong University, Shanghai, China

{cheng.chang, yang\_runzhe, chenlusz, owenzx, kai.yu}@sjtu.edu.cn

## A Figure 1: Example of Successful Dialogue

TASK: ask for <i>italian</i> restaurant in <i>north</i> area & request its <i>phone number</i>					
	Dialogue Turn	Score	$Q^{turn}$	$Q^{succ}$	FP
[1]	System [SLU] welcomemsg()				
	User [Top ASR] Italian food in the north part of town.	0.30	-4.54	27.44	False
[2]	System [SLU] expl-conf(food="italian")				
	User [Top ASR] Yes.	0.99	-2.24	29.09	False
[3]	System [SLU] offer(name="caffe uno") inform(food="italian") inform(area="north")				
	User [Top ASR] The phone number.	0.92	-2.00	28.27	False
[4]	System [SLU] offer(name="caffe uno") inform(food="italian") inform(area="north") Inform(phone="01223314954")				
	User [Top ASR] Does it serve danish italian food.	0.53	-2.41	28.20	False
[5]	System [SLU] offer(name="caffe uno") inform(food="italian") inform(area="north")				
	User [Top ASR] Goodbye.	0.58	0.05	27.42	False

Figure 1: An example of successful dialogue while training without teaching.

## B Algorithm 1: the details of FPT heuristic

## C Figure 2-4: On-line learning process under different teaching schemes

\* Both authors contributed equally to this work.

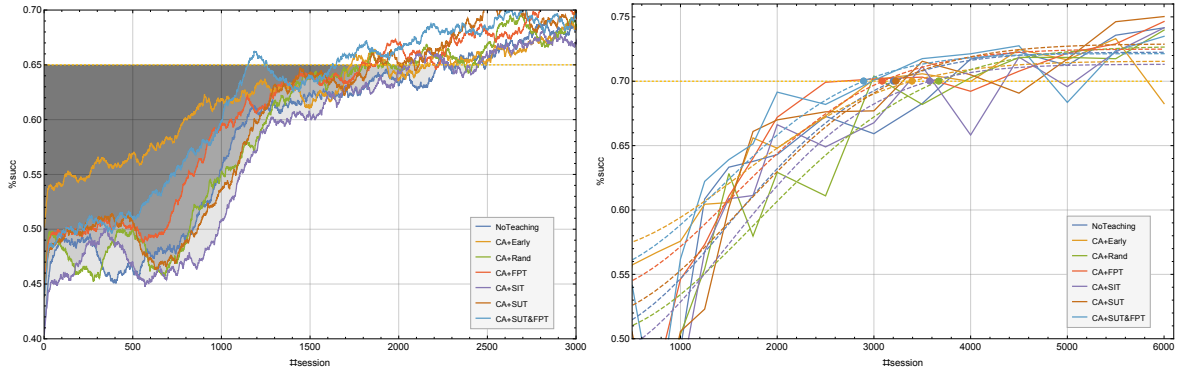


Figure 2: **Left:** On-line learning process under different teaching schemes (CA + different heuristic-s). **Right:** Test curves and fitted empirical learning curves of learning process with different teaching schemes (CA+different heuristic).

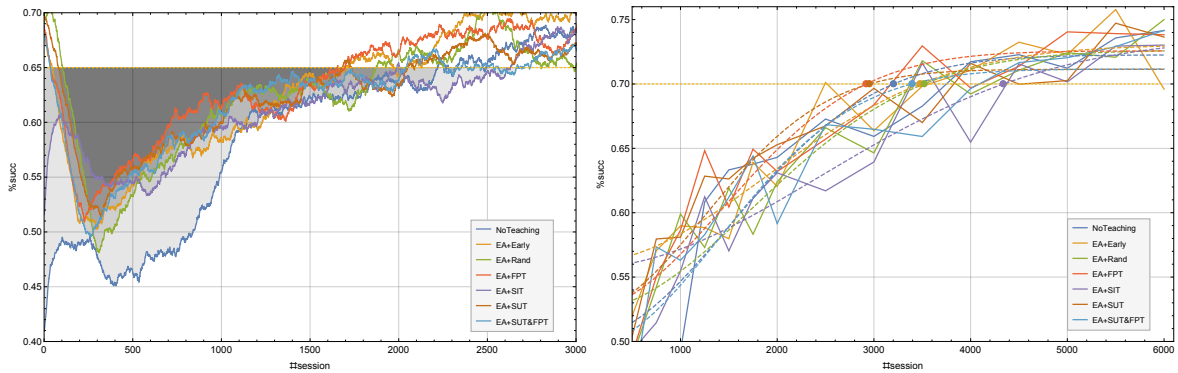


Figure 3: **Left:** On-line learning process under different teaching schemes (EA + different heuristic-s). **Right:** Test curves and fitted empirical learning curves of learning process with different teaching schemes (EA+different heuristic).

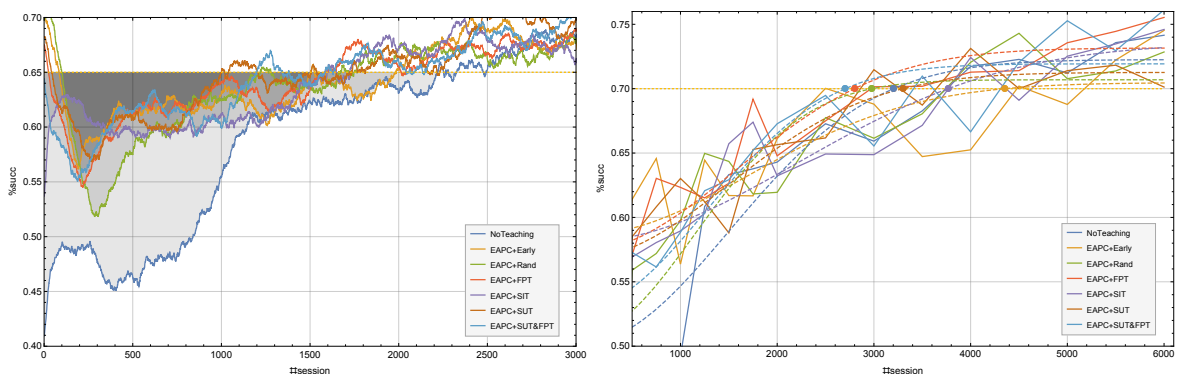


Figure 4: **Left:** On-line learning process under different teaching schemes (EAPC + different heuristic-s). **Right:** Test curves and fitted empirical learning curves of learning process with different teaching schemes (EAPC+different heuristic).

---

**Algorithm 1** Failure Prognosis Based Teaching Heuristic

---

- 1: Initialize replay memory  $\mathcal{D}$
  - 2: Initialize MTL Q-Network,  $Q^{\text{turn}}$  and  $Q^{\text{succ}}$ , with random weights
  - 3: Initialize teaching budget  $c$ , ratio threshold  $\alpha$ , sliding window size  $w$
  - 4: Initialize current teaching strategy (can be any strategy described in section 2.1)
  - 5: Set teaching step  $k \leftarrow 0$
  - 6: **for**  $episode = 1, N$  **do**
  - 7:   Initialize dialogue state  $s_0$
  - 8:   **for**  $t = 0, T$  **do**
  - 9:     Select  $a_t$  randomly with probability  $\epsilon$ , otherwise select:  
     $\text{argmax}_a(Q^{\text{turn}}(s_t, a) + Q^{\text{succ}}(s_t, a))$
  - 10:     **if**  $k < c$  and *failure prognosis is true* according to equation 5 **then**
  - 11:       Ask teacher for advice action  $a_t^{\text{tea}}$
  - 12:        $k \leftarrow k + 1$
  - 13:     **end if**
  - 14:     Update  $a_t$  by current teaching strategy
  - 15:     Take action  $a_t$ , observe  $r_t^{\text{turn}}$  and  $r_t^{\text{succ}}$ , transit to next state  $s_{t+1}$
  - 16:     Update  $r_t^{\text{turn}}, r_t^{\text{succ}}$  according to current teaching strategy
  - 17:     Store  $(s_t, a_t, r_t^{\text{turn}}, r_t^{\text{succ}}, s_{t+1})$  in  $\mathcal{D}$
  - 18:     Sample minibatch of transitions  $e \leftarrow (s_j, a_j, r_j^{\text{turn}}, r_j^{\text{succ}}, s_{j+1})$  from  $\mathcal{D}$
  - 19:     Update  $Q_e^{\text{turn}}$  and  $Q_e^{\text{succ}}$  according to equation 4, with respect to corresponding parameters
  - 20:     Optimize  $(Q_e^{\text{turn}} - Q^{\text{turn}}(s, a; \theta^{\text{turn}}))^2$  and  $(Q_e^{\text{succ}} - Q^{\text{succ}}(s, a; \theta^{\text{succ}}))^2$  simultaneously under MTL structure, using gradient descent.
  - 21:   **end for**
  - 22: **end for**
- 

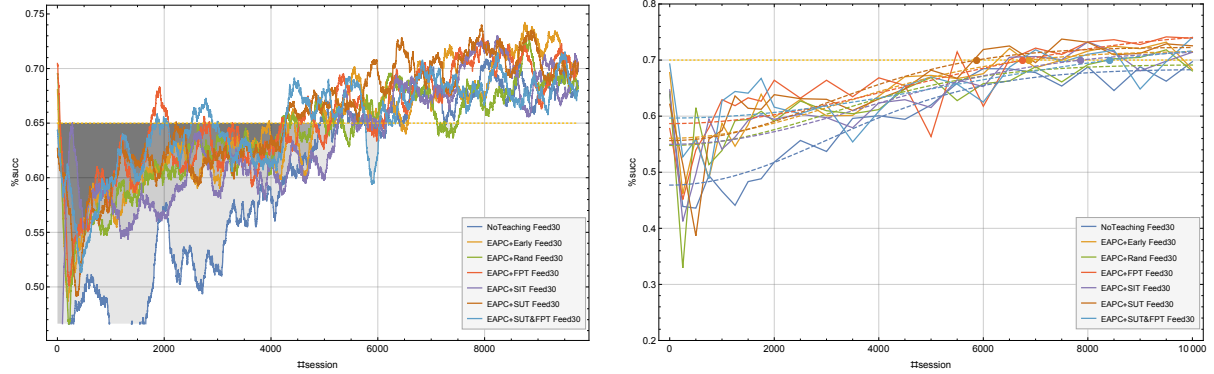
**D Figure 5: On-line learning process with sparse user feedback**

Figure 5: **Left:** On-line learning process under different teaching schemes (EAPC + different heuristics). **Right:** Test curves and fitted empirical learning curves of learning process with different teaching schemes (EAPC+different heuristic). User feedback rate is 30%.