

Autonomous Sub-domain Modeling for Dialogue Policy with Hierarchical Deep Reinforcement Learning

Giovanni Yoko Kristianto¹, Huiwen Zhang²³, Bin Tong¹, Makoto Iwayama¹,
Yoshiyuki Kobayashi¹

¹Hitachi Central Research Laboratory, Tokyo, Japan

²Shenyang Institute of Automation, Shenyang, China

³University of Chinese Academy of Sciences, Beijing, China

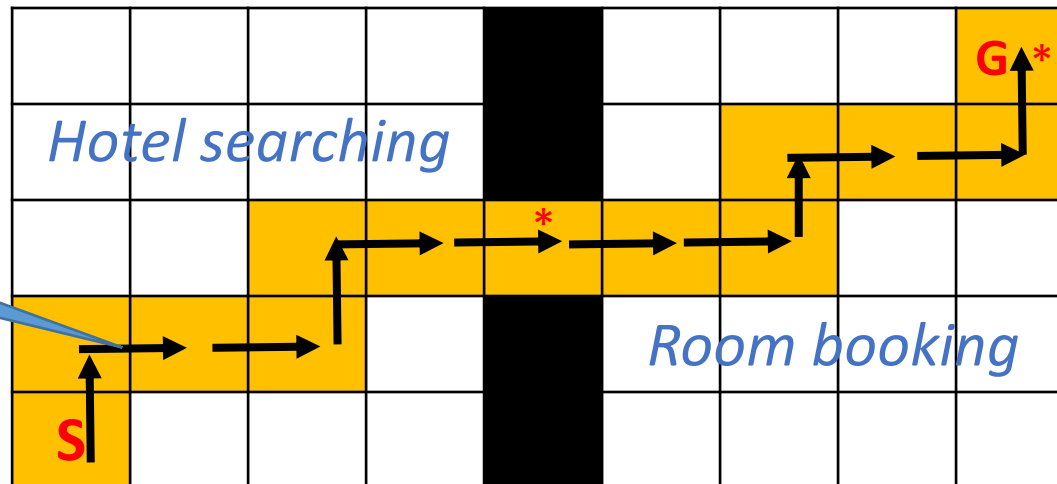
Introduction: Motivation

- Dialogue system for composite task-oriented domains
 - Hotel reservation: hotel searching + room booking
 - Travel planning: flight reservation + hotel reservation + car renting
- Challenges in reinforcement learning for composite domain
 - Composite dialogues have larger state-actions space
 - The state trajectory is longer (i.e. need more turns)
 - Sparser rewards

Introduction: Solution

- To decompose the composite domain into multiple sub-domains
e.g. hotel reservation -> hotel searching + room booking
- To provide internal reward function for each sub-domain, so that the rewards are less sparse

S: Do you need parking space?
U: Yes

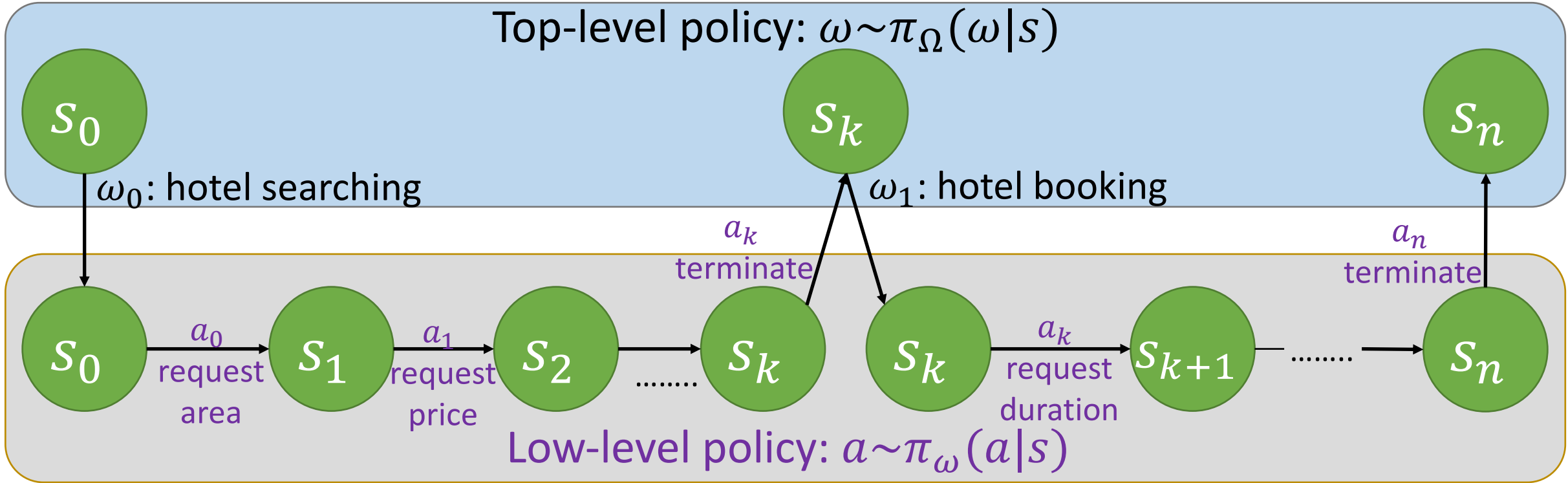


S: starting state

G: goal state

*: subdomain's goal state

Dialogue Policy Execution using Sub-domains



Previous Works

- Peng et al. 2017. *Composite Task-Completion Dialogue Policy Learning via Hierarchical Deep RL*. In EMNLP '17
- Budzianowski et al. 2017. *Sub-domain Modelling for Dialogue Management with HRL*. In SIGDIAL '17

Summary:

- Manually defining composing sub-domains (and the internal reward functions)

Our Approach

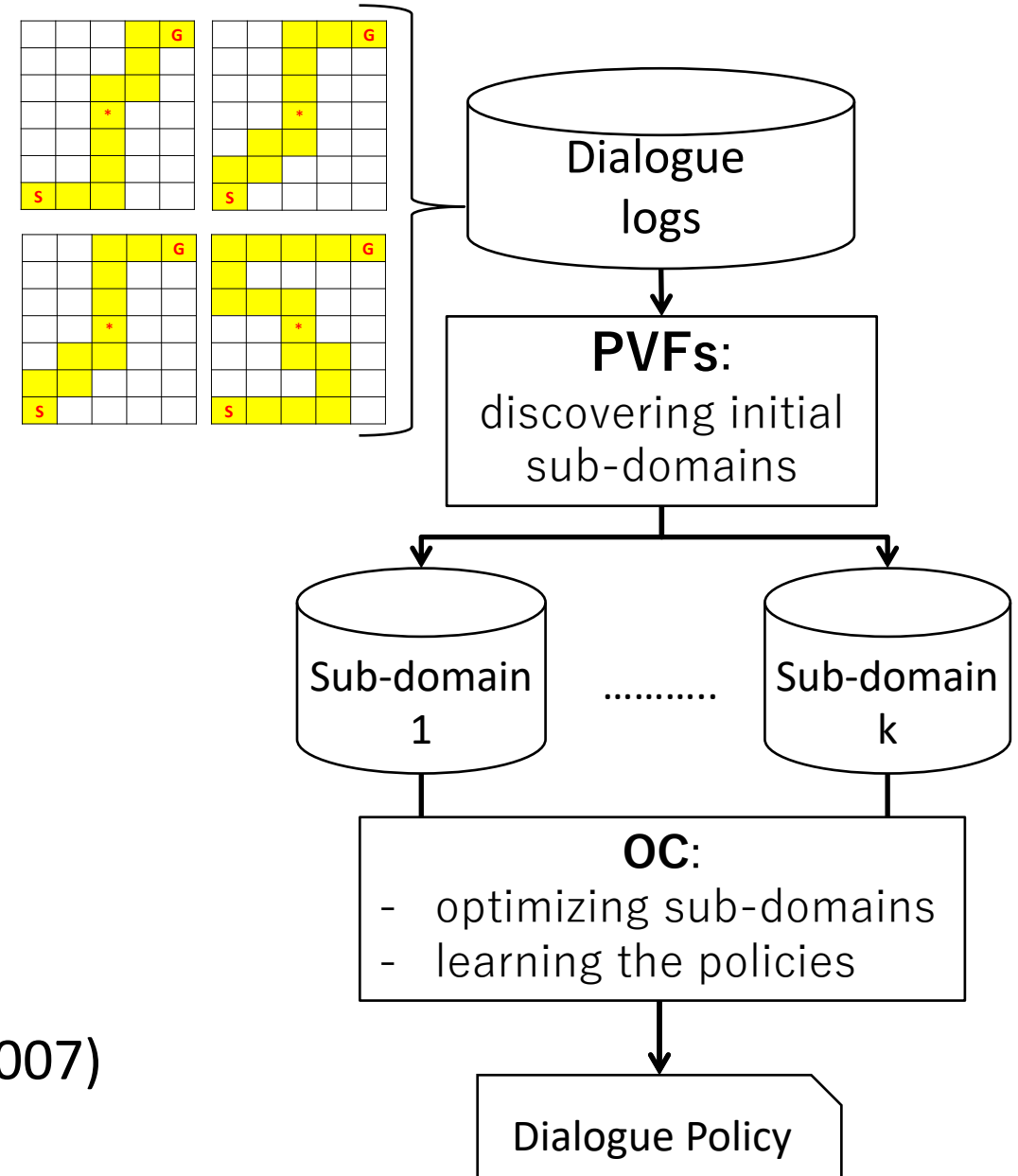
- To autonomously discover sub-domains (and their internal reward functions), and
- To utilize them for training dialogue policy

Input:

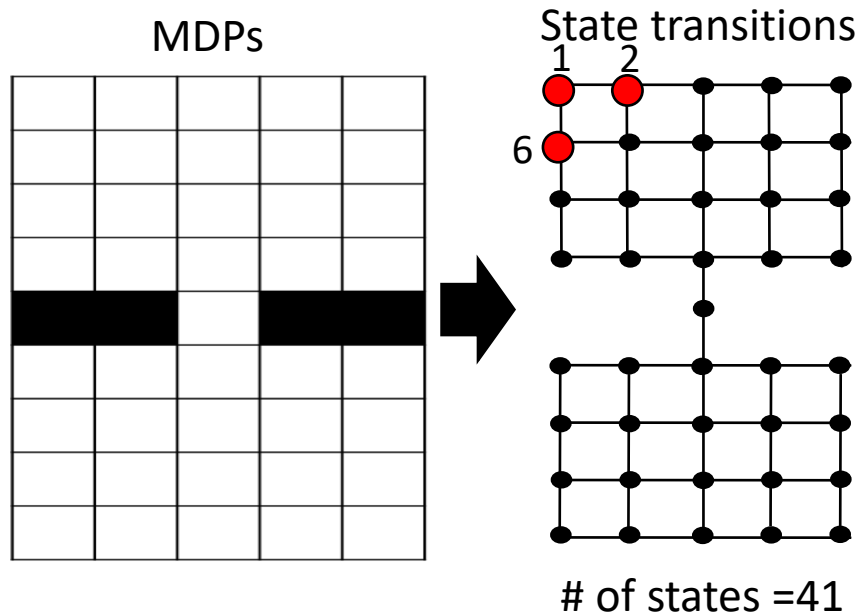
- Dialogue logs
- k: number of sub-domains

Methods:

- Proto-value functions (PVFs; Mahadevan, 2007)
- Option Critic (OC; Bacon, 2017)



Proto-value function (PVFs) (Mahadevan, 07; Machado, 17)

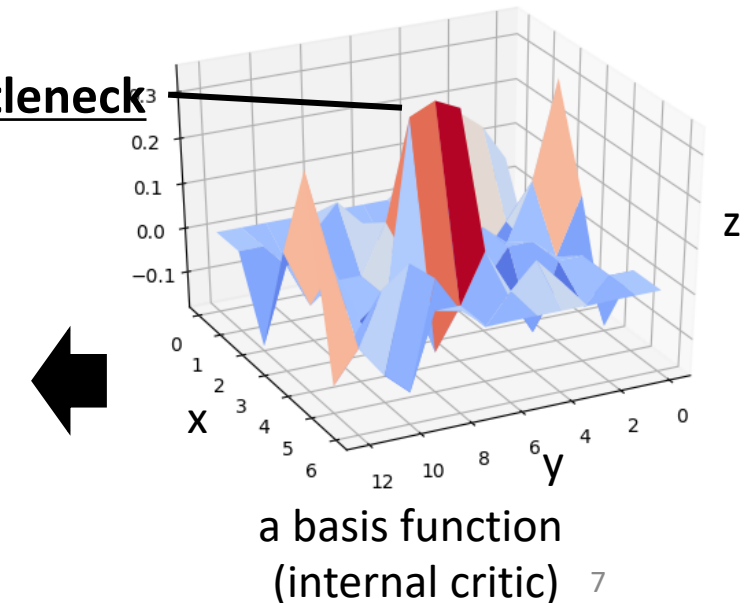
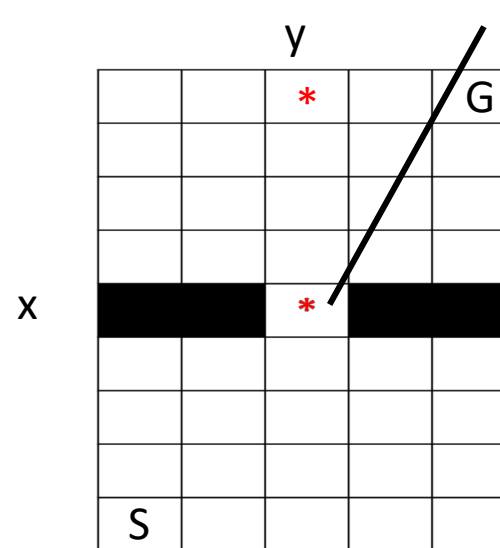


Adjacency matrix (A)

	1	2	3	4	5	6	..	41
1	0	1	0	0	0	1		
2	1	0	1	0	0	0	...	
3	0	1	0	1	0	0		
		\vdots			\vdots		\vdots	
39	0	0	0	0	0	0		
40	0	0	0	0	0	0	...	
41	0	0	0	0	0	0		

1. Get Laplacian Matrix
 $L = D - A$
2. Obtain basis functions from L by eigen-decomposition

- PVFs use no external reward
- Discovered sub-goals are task-independent
 - *But, we need to maximize external rewards specific to our task!*

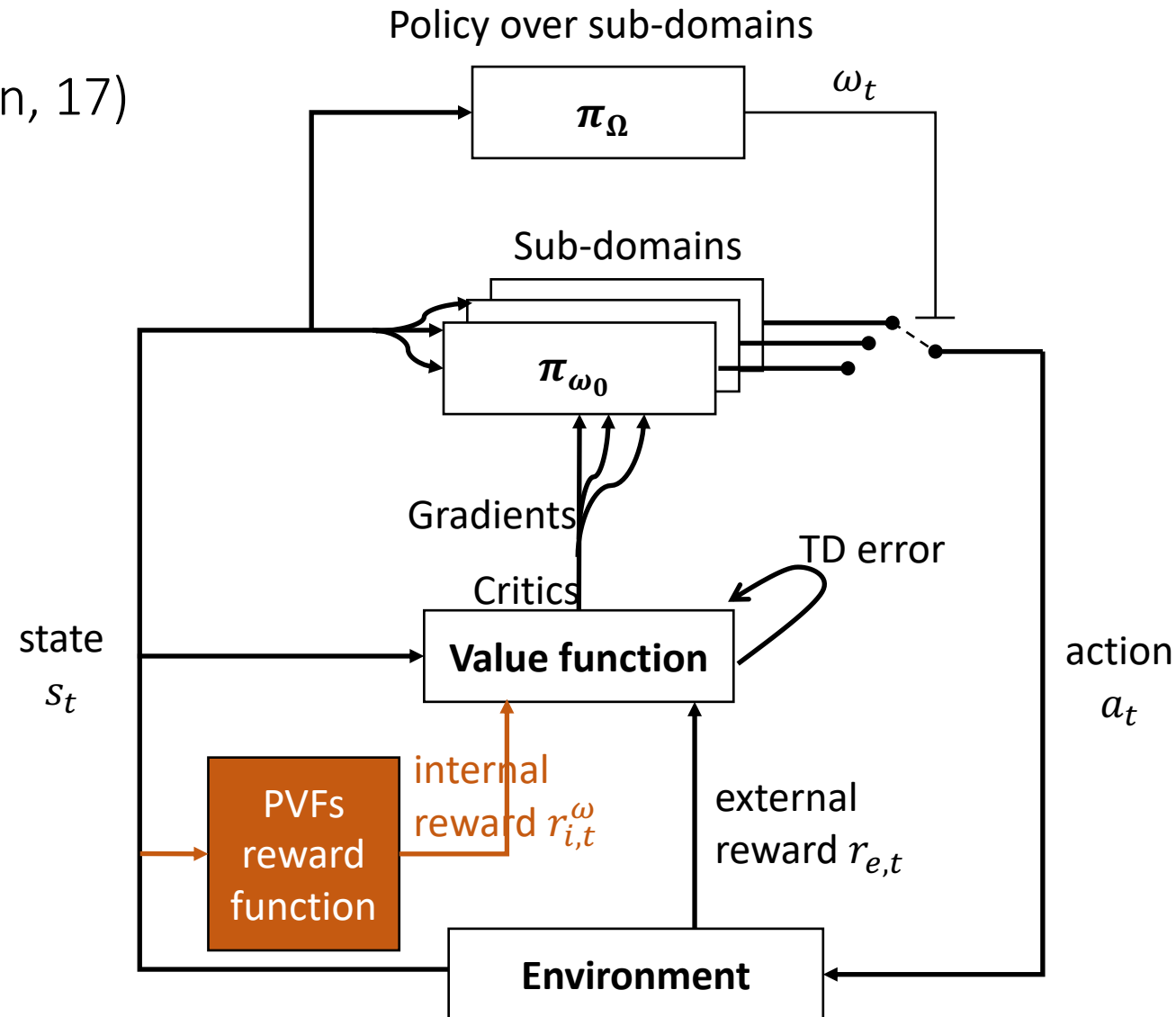


Our Approach: Option-Critic (OC) (Bacon, 17)

- OC learns sub-domains using external reward only

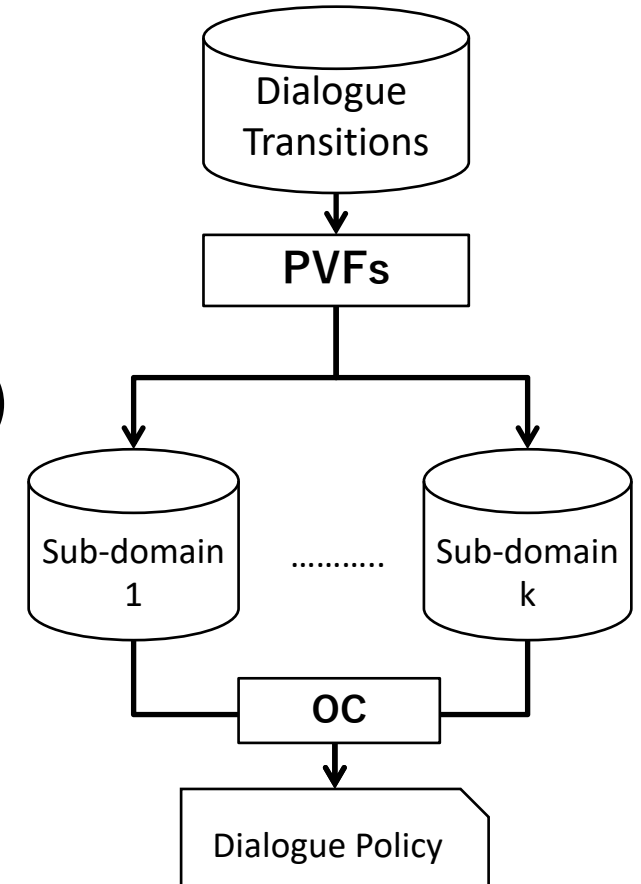
Combining PVFs and OC:

- To update the critics by $\alpha r_{i,t}^\omega + (1 - \alpha)r_{e,t}$



Experiment

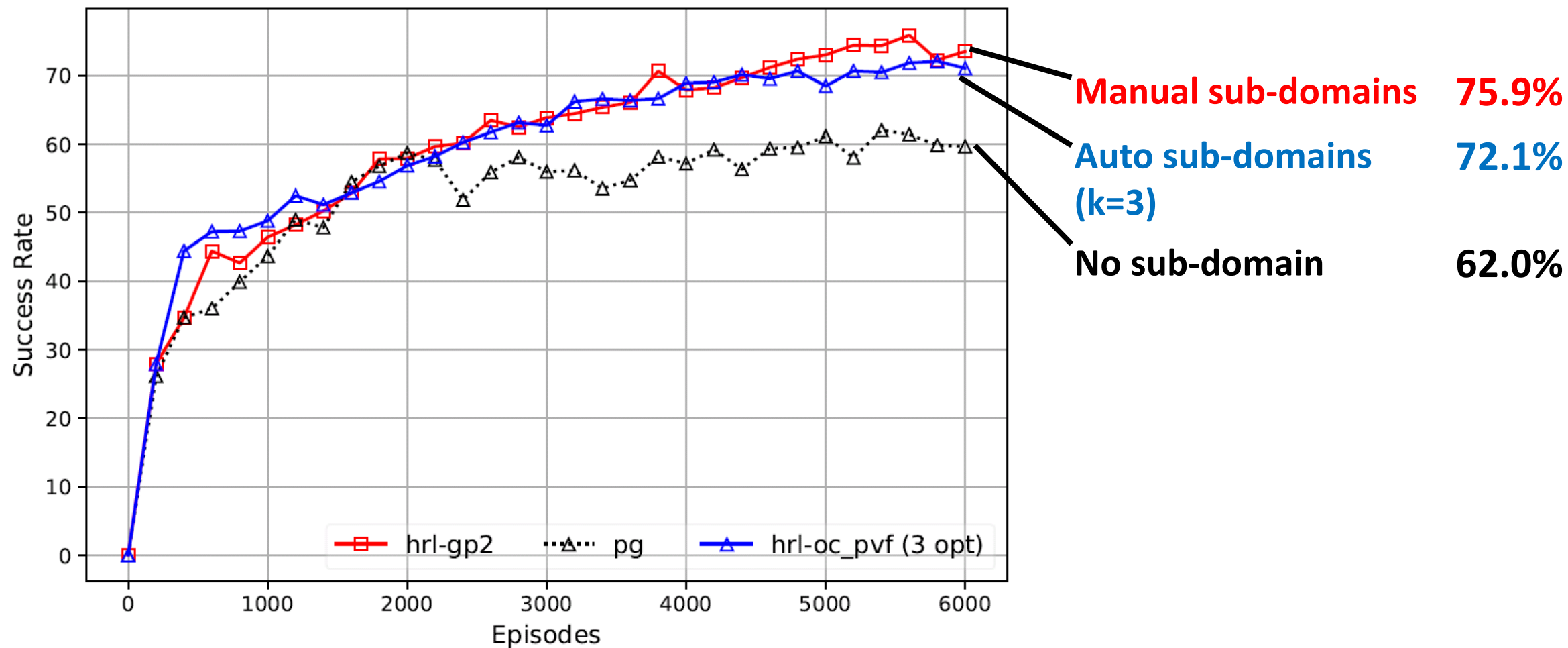
- Impact on success rate
 - Policy without sub-domain (Reinforce)
 - Policy with manual sub-domain (Gaussian Process RL)
 - Policy with autonomous sub-domains (proposed approach)
- Qualitative analysis of the discovered sub-domains



Experiment: Setup

- Dialogue domain: Hotel Reservation
 - # of constraint slots: 13
 - # of system actions: 44
- Training and testing utilized simulator
 - Behavior: searching for a hotel, followed by booking and/or payment
- Implementation
 - Sample transitions used by PVFs: 1,000 (sampled from 100 dialogue samples)

Results: Success Rate



Results: goals of discovered sub-domains are intuitive

<p>.... U: request(address) S: inform(address=xxx) U: hello(day=..., duration=..., name=..., task=booking)</p> <p>user ends 'searching', then enters 'booking'</p>	<p>.... U: inform(address=..., cardNo=..., name=..., task=payment) S: inform(paid=true) U: bye()</p> <p>user ends 'payment'</p>
<p>.... U: hello(day=..., duration=..., name=..., task=booking) S: inform() U: request_alternative()</p> <p>user ask alternative information</p>	<p>.... U: inform(peopleNo=..., task=booking) S: inform(booked=true) U: bye()</p> <p>user ends 'booking'</p>

Conclusion

- We proposed a dialogue policy learning framework consisting of proto-value functions and option-critic
 - It performs comparable with using manually modeled sub-domains
 - The discovered sub-domains are intuitive
- Future works
 - Human evaluation
 - Transfer learning by exploiting sub-domains