



“Think Before You Speak”: Improving Multi-Action Dialog Policy by Planning Single-Action Dialogs

Shuo Zhang¹, Junzhou Zhao^{1*}, Pinghui Wang^{1*}, Yu Li¹, Yi Huang², Junlan Feng²

¹MOE KLINNS Lab, Xi'an Jiaotong University, Xi'an 710049, P. R. China

²JIUTIAN Team, China Mobile Research, Beijing 100053, P. R. China

{zs412082986, liyu1998}@stu.xjtu.edu.cn, {junzhou.zhao, phwang}@mail.xjtu.edu.cn,
{huangyi, fengjunlan}@chinamobile.com

Code: <https://github.com/ShuoZhangXJTU/PEDP>.

(IJCAI-2022)

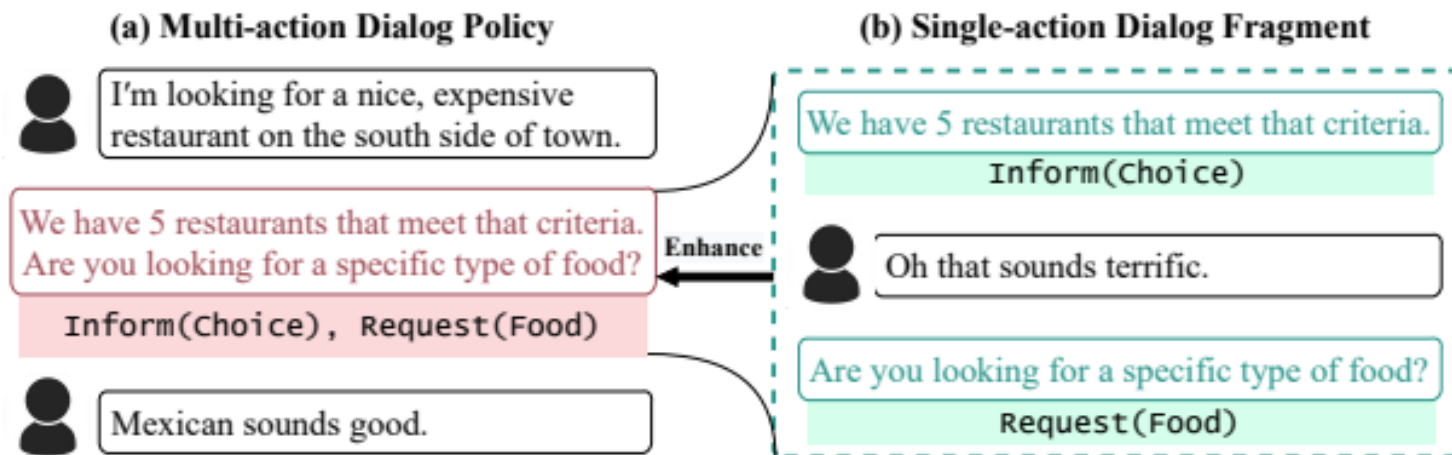




1. Introduction
2. Approach
3. Experiments



Introduction



We propose Planning Enhanced Dialog Policy (PEDP), a novel multi-task learning framework that learns single action dialog dynamics to enhance multi-action prediction.

Figure 1: (a) Example dialog under multi-action dialog policy². We propose to learn single-action dialog dynamics (b) to model conditional act combination patterns and enhance multi-action prediction.

²A dialog policy responses by predicting atomic dialog actions represented as “Domain-Intent(Slot)” phrases. We omit the domain (“restaurant”) for clarity.

- 一个宏动作，它是一组独立的原子对话动作，用作当前系统响应。
- 每个原子对话动作都是域名、动作类型和插槽名称的串联，例如“hotel-inform-area”。

Approach

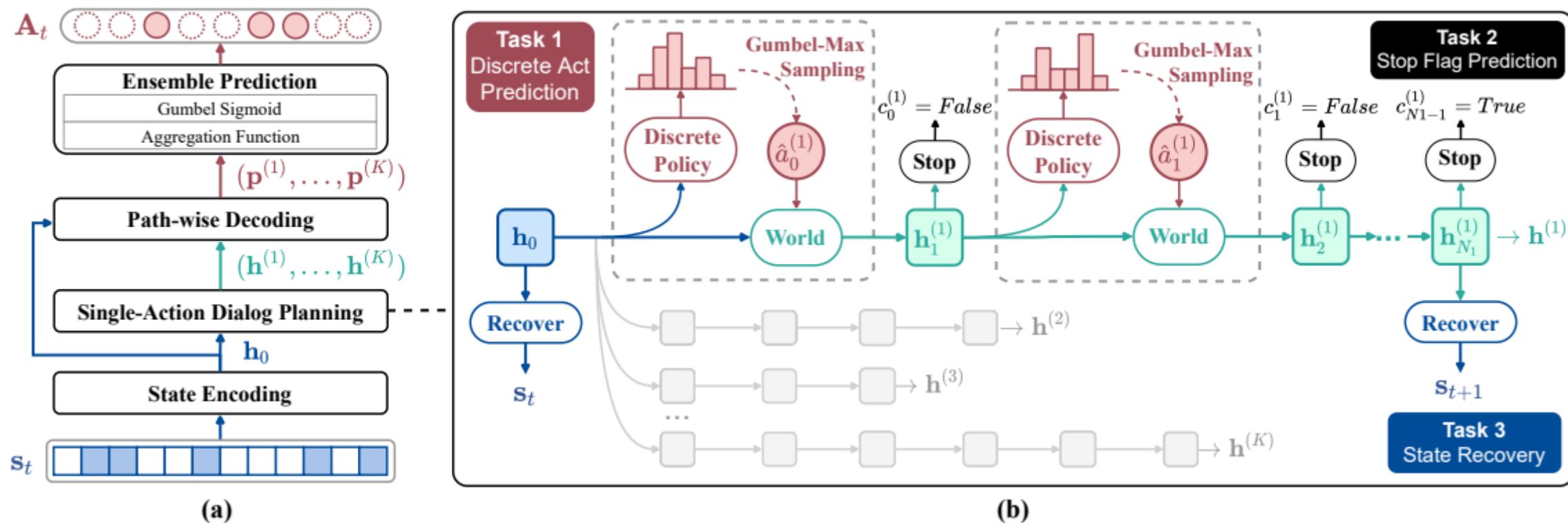
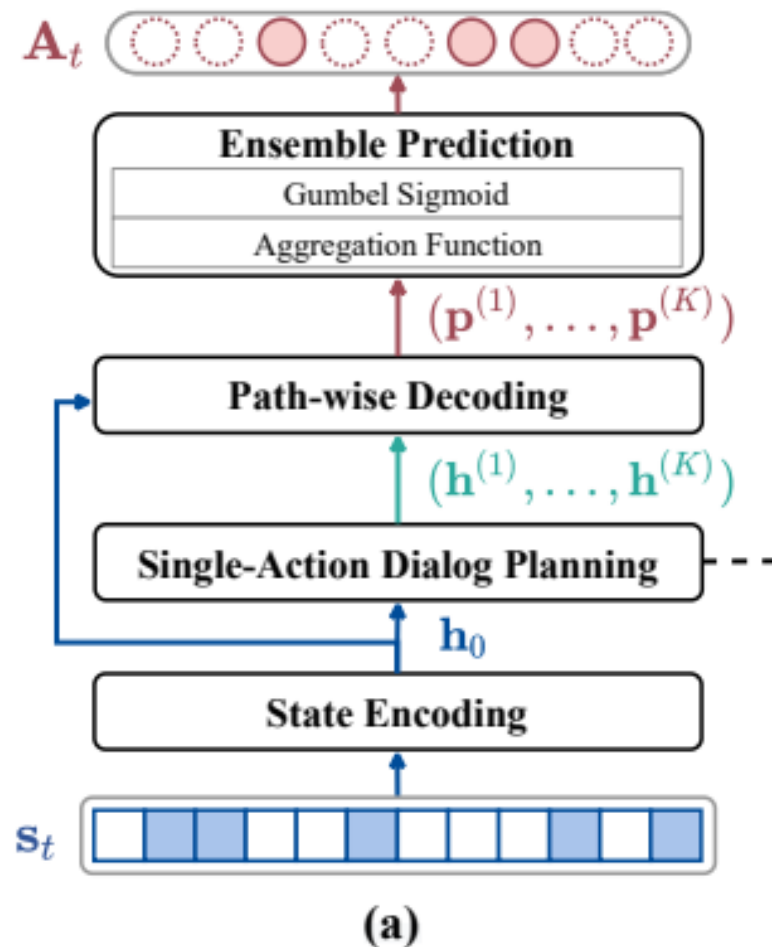


Figure 2: (a) The Planning-Enhanced Dialog Policy (PEDP) framework. It utilizes a *single action dialog planning* module (b) to incorporate contextually relevant contents before multi-action prediction. A total of K single-action dialog procedures are planned, with the k -th path looking ahead N_k steps under single-action dialog dynamics. At each step, the discrete policy model predicts an atomic dialog action a_n given the previous dialog state embedding h_{n-1} . The world model, which simulates user behavior, responds to the predicted action a_n and updates the dialog state embedding from h_{n-1} to h_n .

Approach



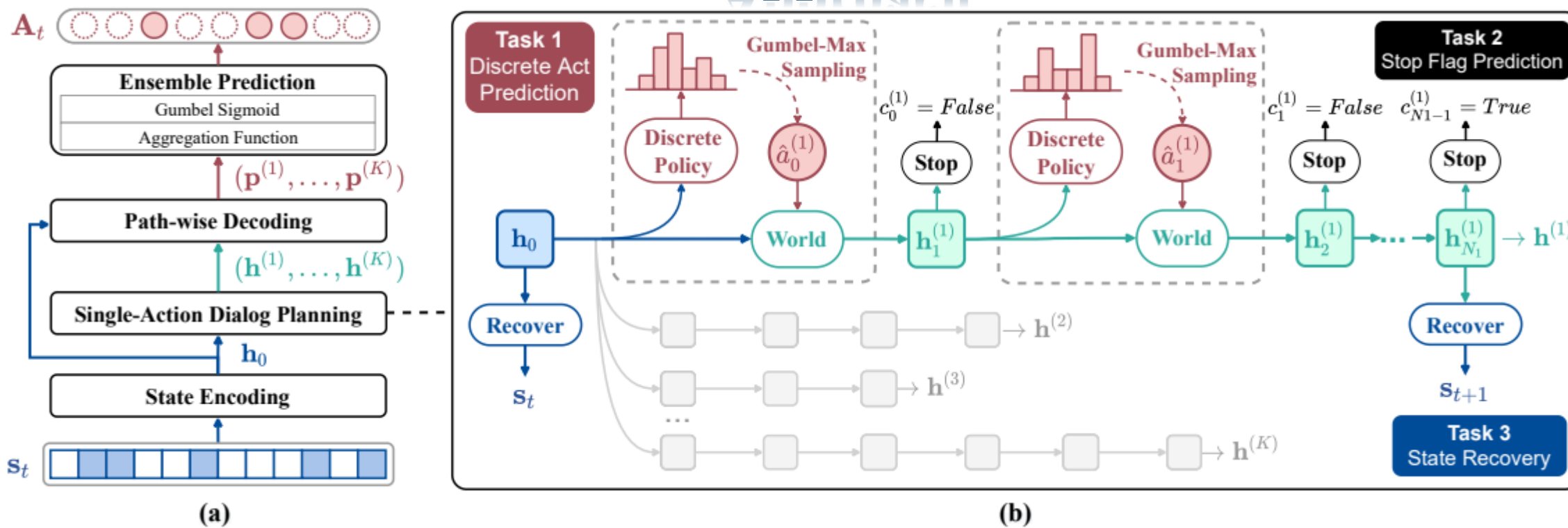
\mathbf{s}_t into a dialog state embedding \mathbf{h}_t . Given the current state embedding \mathbf{h}_t , we plan K independent single-action dialog paths, and the k -th dialog path is represented by a vector $\mathbf{h}^{(k)}$, $k = 1, \dots, K$. Our model then decodes each dialog path to a probability distribution over atomic dialog actions, i.e., $\mathbf{p}^{(k)}$. Finally, these distributions $\{\mathbf{p}^{(k)}\}_{k=1}^K$ are aggregated to form a unified distribution, from which atomic dialog actions in the macro-action \mathbf{A}_t are sampled.

State Encoding

$$\mathbf{h}_t = \text{FFN}_{enc}(\mathbf{s}_t) = \text{ReLU}(\mathbf{s}_t W_1 + b_1) W_2 + b_2. \quad (1)$$

In what follows, this dialog state embedding \mathbf{h}_t will serve as the initial dialog state embedding for planning.

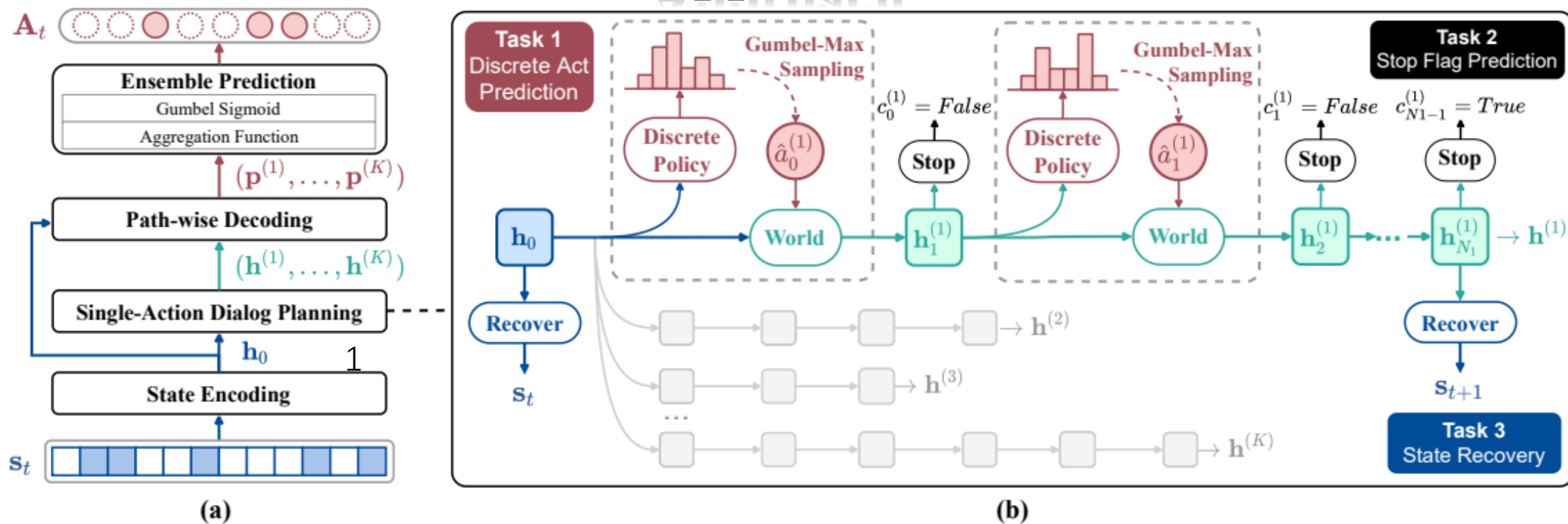
Approach



Single-Action Dialog Planning

look ahead several steps. Let $h_{t,n}^{(k)}$ denote the dialog state embedding at the n -th step in the k -th dialog for $n = 0, \dots, N_k$ where N_k is the length of the k -th dialog, and $h_{t,0}^{(k)} = h_t$, $h_{t,N_k}^{(k)} = h^{(k)}$. The last dialog state embedding $h^{(k)}$ estimates the hidden vector of the future dialog state s_{t+1} and summarizes the planned single-action dialog. In what follows, we describe how to plan a single step from $h_{t,n}^{(k)}$ to obtain $h_{t,n+1}^{(k)}$.

Approach

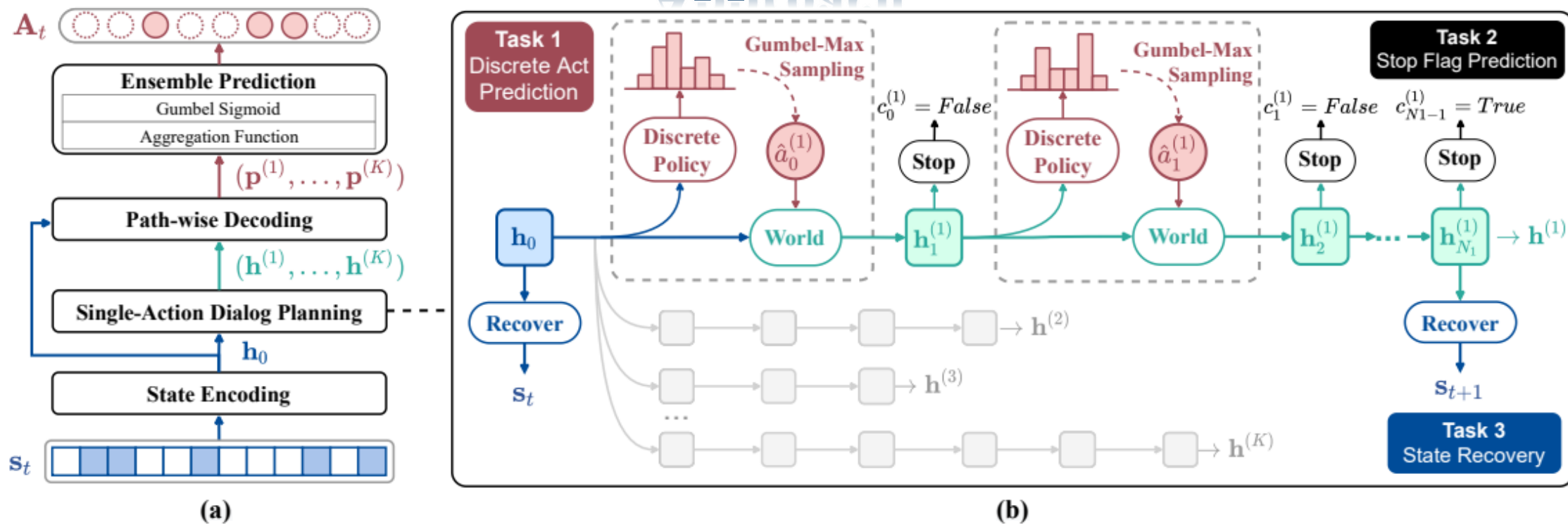


Single-Action Dialog Planning

$$\begin{aligned}
 a_n &= \text{DP}(\mathbf{h}_n) \triangleq \text{GumbelSoftmax}^{(\tau_d)}(\mathbf{h}_n W_d + b_d) \\
 \mathbf{h}_{n+1} &= \text{World}(\mathbf{h}_n, a_n) \triangleq \text{GRU}(\mathbf{h}_n, \text{Emb}(a_n)). \quad (2)
 \end{aligned}$$

Here, DP is implemented as a single linear layer followed by a Gumbel-Softmax function [Jang *et al.*, 2016] parameterized by τ_d . The Gumbel-Softmax function draws an atomic dialog action sample from a categorical distribution, diversifying the planned dialogs. τ_d is selected to balance the approximation bias and the magnitude of gradient variance. The world model is implemented using a GRU to model dialog state transitions, and $\text{Emb}(a_n)$ denotes the embedding vector of atomic dialog action a_n .

Approach



Single-Action Dialog Planning

$$c_n = \text{GumbelSoftmax}^{(\tau_s)}(\text{FFN}_{s_t}([\mathbf{h}_0 : \mathbf{h}_{n+1}]))$$

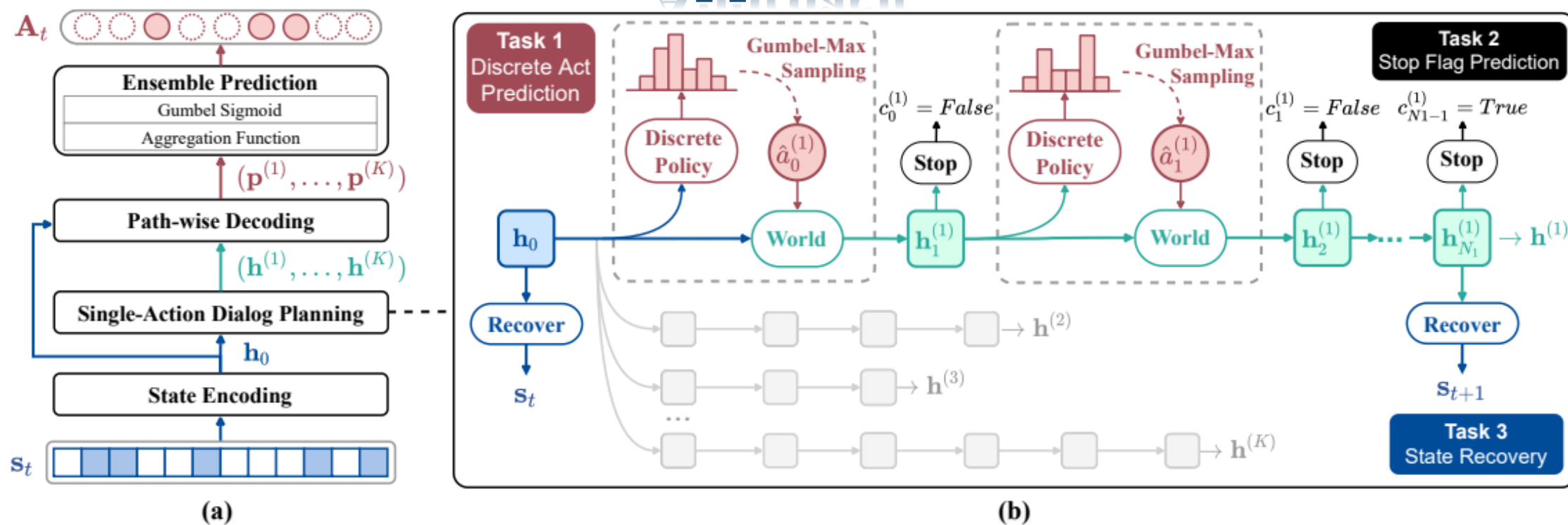
where c_n is a binary variable, “:” denotes vector concatenation, and FFN is a 2-layer fully-connected feed-forward network using the ReLU activation function in the middle layer.

$$\mathbf{s}_t = \text{Recover}(\mathbf{h}_0) \quad (3)$$

$$\mathbf{s}_{t+1} = \text{Recover}(\mathbf{h}_N)$$

Here, Recover is implemented by a 2-layer FFN and is only used during the training stage.

Approach

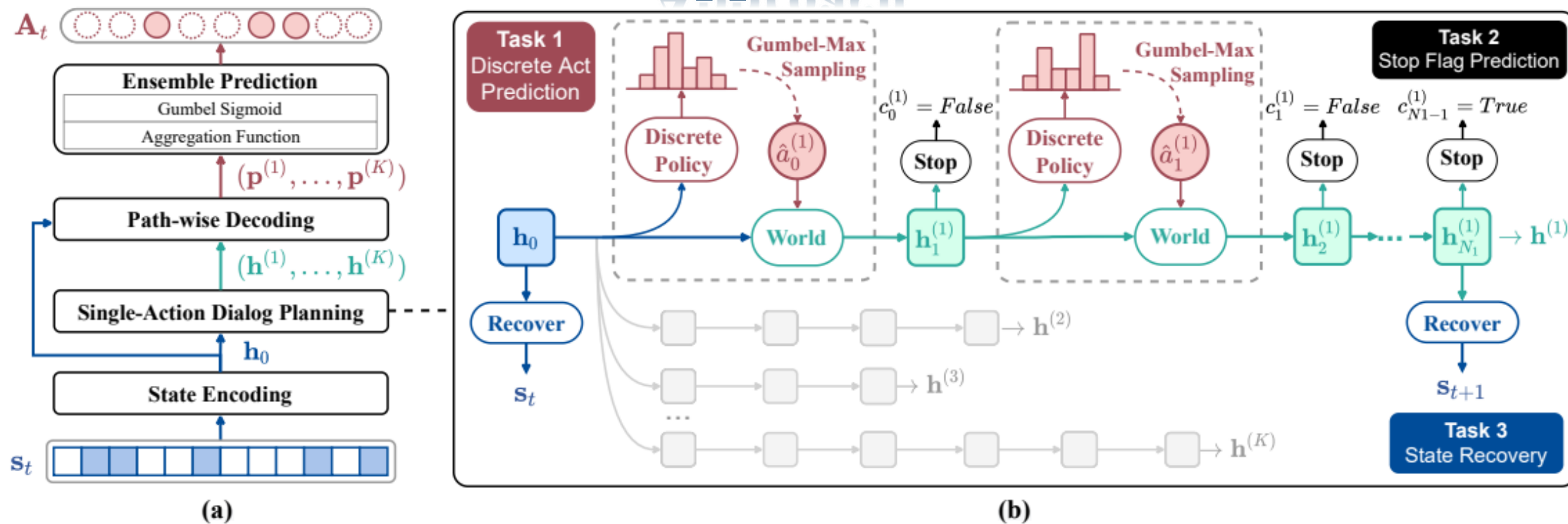


Path-wise Decoding

Specifically, we instantiate the decoder $\mathbf{p}^{(k)} = [\mathbf{p}_1^{(k)} : \dots : \mathbf{p}_M^{(k)}]$, where k refers to the planned path and M is the size of the action space. Each $\mathbf{p}_m^{(k)}, m = 1, \dots, M$ is a vector computed as:

$$\mathbf{p}_m^{(k)} = \text{FFN}_m^{\text{dec}}([\mathbf{h}_0 : \mathbf{h}^{(k)}]).$$

Approach



Ensemble Prediction

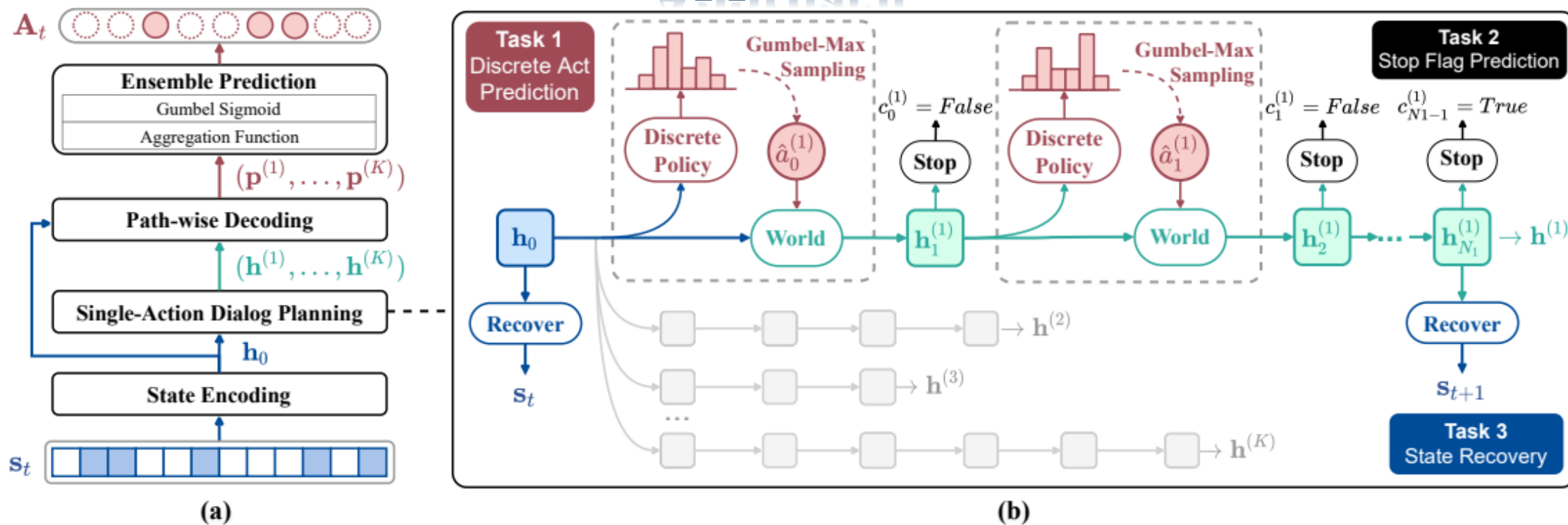
$$\mathbf{P}_t = \text{Aggr}(\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(K)})$$

where $\text{Aggr}(\cdot)$ is the mean average in our case.

$$\mathbf{A}_t = \text{GumbelSigmoid}(\mathbf{P}_t) = \frac{e^{(\mathbf{P}_t + g_1)/\tau}}{e^{(\mathbf{P}_t + g_1)/\tau} + e^{(\mathbf{P}_t + g_2)/\tau}}$$

Here $\text{GumbelSigmoid}(\cdot)$ is a modification of the Gumbel-Softmax function, regarding sigmoid as a softmax with two logits p and 0 . τ denotes the temperature factor, g_1 and g_2 are Gumbel noises.

Approach



Training

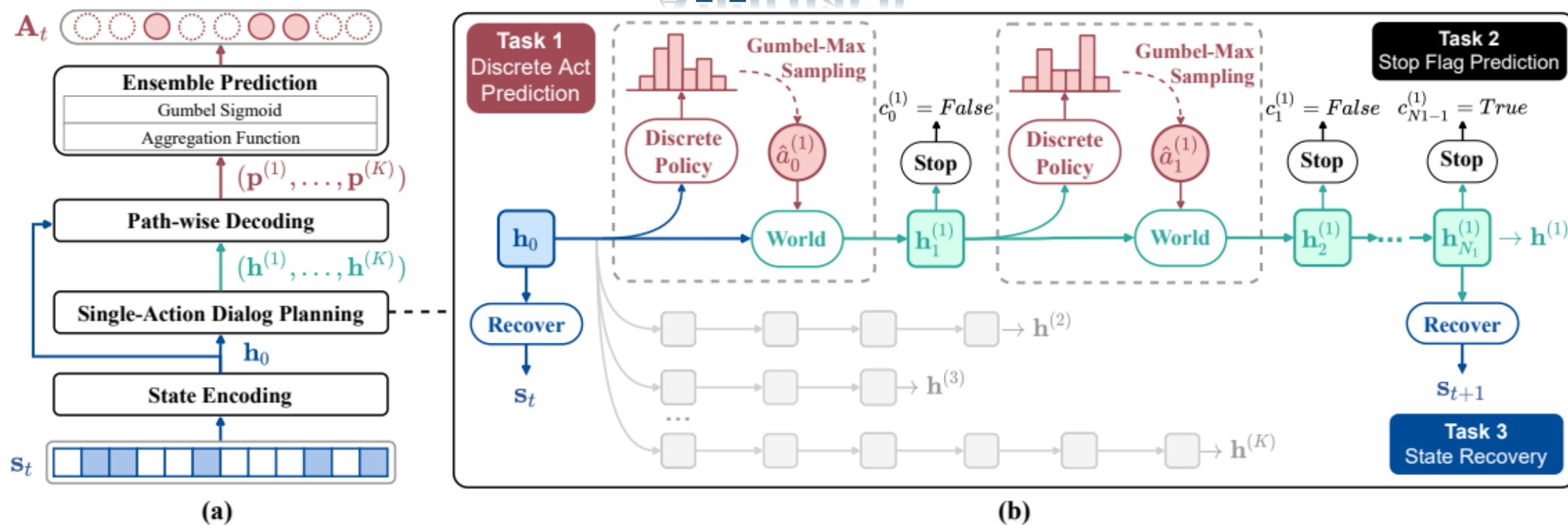
Task 1: Discrete Act Prediction (DAP)

$$p(\mathbf{a}|\mathbf{h}_0) = p_\theta(a_0|\mathbf{h}_0) \prod_{n=1}^{N-1} \underbrace{p_\theta(a_n|\mathbf{h}_n)}_{\text{DAP}} \underbrace{p_\phi(\mathbf{h}_n|a_{n-1}, \mathbf{h}_{n-1})}_{\text{state transition}}$$

$$\mathbf{a} = (a_0, \dots, a_{N-1})$$

where θ and ϕ denotes trainable parameters for the discrete dialog policy model and the world model, respectively.

Approach



Training

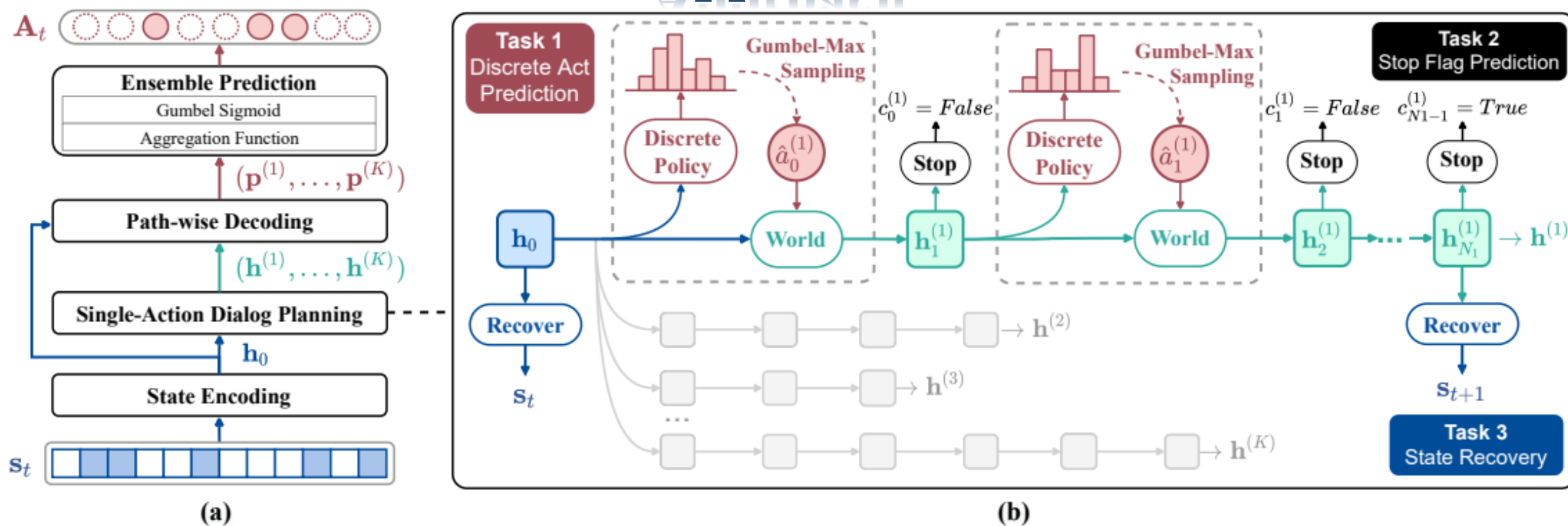
Task 2: Stop Flag Prediction (SFP)

$$p(\mathbf{c}|\mathbf{h}_0) = \prod_{n=0}^{N-1} \underbrace{p_\gamma(c_n|\mathbf{h}_{n+1}, \mathbf{h}_0)}_{\text{SFP}} \underbrace{p_{\phi, \theta}(\mathbf{h}_{n+1}|\mathbf{h}_n)}_{\text{1-step planning}}$$

$$\mathbf{c} = (c_0, \dots, c_{N-1})$$

where γ parameterizes the stop prediction model, the joint probability of $p_{\phi, \theta}(\mathbf{h}_{n+1}|\mathbf{h}_n)$ is factorized as $p_\phi(\mathbf{h}_{n+1}|a_n, \mathbf{h}_n)p_\theta(a_n|\mathbf{h}_n)$ of state transition and discrete act prediction.

Approach



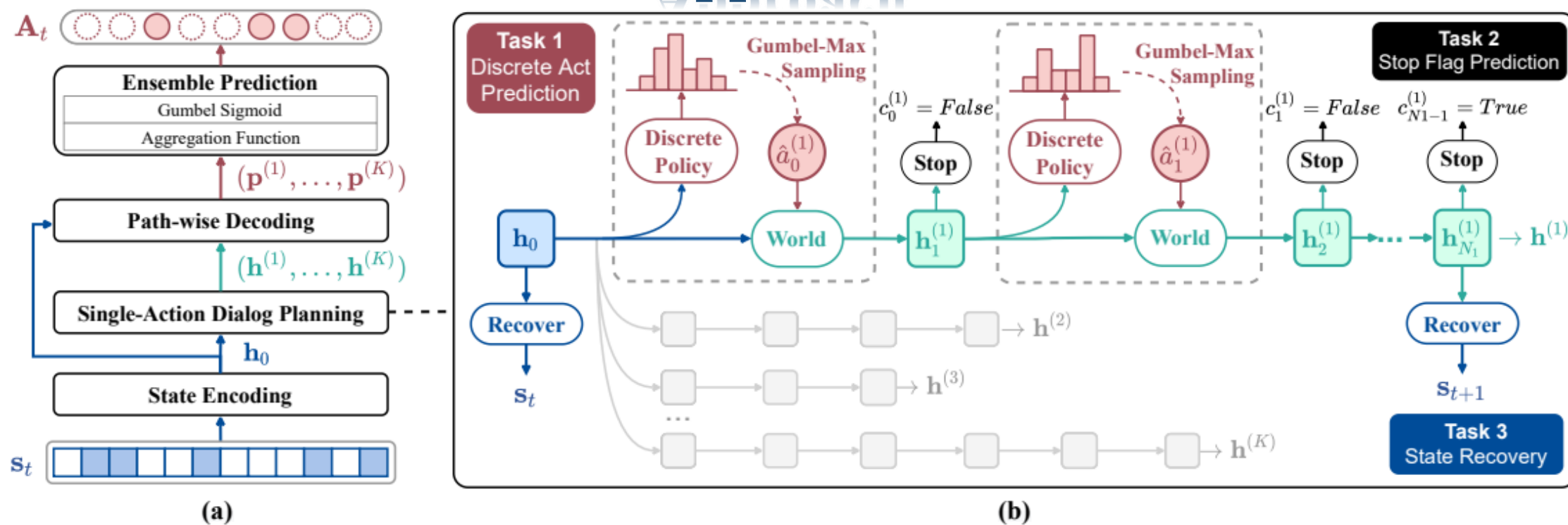
Training Task 3: State Recovery (SR)

$$p(s_t) = \underbrace{p_{\zeta}(s_t | h_0)}_{\text{SR}} \underbrace{p_{\eta}(h_0 | s_t)}_{\text{state encoding}}$$

$$p(s_{t+1} | s_t) = \underbrace{p_{\zeta}(s_{t+1} | h_N)}_{\text{SR}} \underbrace{p_{\eta}(h_0 | s_t)}_{\text{state encoding}} \prod_{n=0}^{N-1} \underbrace{p_{\phi, \theta}(h_{n+1} | h_n)}_{\text{1-step planning}}$$

where η and ζ denotes trainable parameters for state encoder and the Recover, respectively. The joint probability $p_{\phi, \theta}(h_{n+1} | h_n)$ is the same as explained in Task 2.

Approach



Training

Task 4: Multi-Action Prediction (MAP)

$$p(\mathbf{A}_t | s_t) = \underbrace{p_\omega(\mathbf{A}_t | \mathbf{h}_0, \mathbf{h}_N)}_{\text{MAP}} \underbrace{p_\eta(\mathbf{h}_0 | s_t)}_{\text{state encoding}} \prod_{n=0}^{N-1} \underbrace{p_{\phi, \theta}(\mathbf{h}_{n+1} | \mathbf{h}_n)}_{\text{1-step planning}}$$

where ω denotes trainable parameters for the decoder. The rest is the same as explained in Task 3.

Experiments

Agent	MultiWOZ				
	Turn	Match	Rec	F1	Success
DiaMultiClass	11.46 \pm 0.56	0.68 \pm 3.9%	0.81 \pm 3.2%	0.81 \pm 2.1%	67.3 \pm 3.69
+ sample	9.23 \pm 0.2	0.82 \pm 1.1%	0.90 \pm 1.8%	0.77 \pm 1.2%	81.4 \pm 1.78
DiaSeq (beam)	9.06 \pm 0.67	0.81 \pm 0.4%	0.9 \pm 1.2%	0.86 \pm 0.9%	81.4 \pm 0.16
greedy	10.35 \pm 0.04	0.68 \pm 1.5%	0.80 \pm 0.5%	0.77 \pm 0.5%	67.7 \pm 1.02
+ sample	8.82 \pm 0.1	0.86 \pm 0.6%	0.93 \pm 0.4%	0.81 \pm 0.5%	86.9 \pm 0.49
DiaMultiDense	9.66 \pm 0.15	0.85 \pm 0.6%	0.94 \pm 0.4%	0.87 \pm 0.6%	86.3 \pm 0.64
- sample	12.75 \pm 0.77	0.61 \pm 6%	0.72 \pm 5.4%	0.80 \pm 2.3%	58.4 \pm 6.05
gCAS	11.69 \pm 0.53	0.56 \pm 1.4%	0.72 \pm 0.4%	0.76 \pm 1.4%	58.8 \pm 2.82
GP-MBCM ⁵	2.99	0.44	-	0.19	28.9
ACER ⁵	10.49	0.62	-	0.78	50.8
PPO ⁵	15.56	0.60	0.72	0.77	57.4
ALDM ⁵	12.47	0.69	-	0.81	61.2
GDPL	7.54 \pm 0.43	0.84 \pm 0.9%	0.89 \pm 2.2%	0.88 \pm 1.2%	83.2 \pm 1.48
DiaAdv	8.90 \pm 0.18	0.87 \pm 0.9%	0.94 \pm 0.75%	0.85 \pm 0.58%	87.6 \pm 0.9
- sample	11.9 \pm 0.88	0.62 \pm 5.9%	0.73 \pm 4.6%	0.80 \pm 2.1%	61.7 \pm 5.59
PEDP	8.69 \pm 0.15	0.88 \pm 1.3%	0.97 \pm 0.4%	0.87 \pm 1.1%	90.6 \pm 0.68
- planning	9.66 \pm 0.15	0.85 \pm 0.6%	0.94 \pm 0.4%	0.87 \pm 0.6%	86.3 \pm 0.64
- ensemble	9.25 \pm 0.43	0.88 \pm 1.97%	0.96 \pm 0.8%	0.85 \pm 2.5%	89.1 \pm 1.74
- sample	8.85 \pm 0.22	0.82 \pm 2.5%	0.93 \pm 1.4%	0.86 \pm 1.6%	83.4 \pm 1.01

Table 1: Interactive evaluation results. We simulate 1,000 dialogs per run and report the mean and standard deviation over 5 runs.

Experiments

Agent	MultiWOZ			SGD (scaling)		
	F1%	Precision%	Recall%	F1%	Precision%	Recall%
DiaMultiClass	39.41 \pm 1.08	54.59 \pm 1.71	34.32 \pm 1.32	58.09 \pm 0.63	81.29 \pm 1.13	46.29 \pm 0.57
+ sample	38.91 \pm 0.74	47.28 \pm 0.68	37.56 \pm 1.08	58.03 \pm 0.64	81.48 \pm 0.18	46.14 \pm 0.80
DiaSeq (beam)	44.64 \pm 2.08	51.91 \pm 0.99	43.66 \pm 2.27	63.13 \pm 0.18	86.04 \pm 0.5	50.83 \pm 0.30
greedy	48.34 \pm 0.45	54.71 \pm 0.21	48.84 \pm 0.84	63.21 \pm 0.35	86.31 \pm 0.7	50.85 \pm 0.40
+ sample	37.82 \pm 0.45	43.02 \pm 0.48	38.91 \pm 0.64	62.64 \pm 1.03	85.54 \pm 1.62	50.40 \pm 0.76
DiaMultiDense	35.92 \pm 0.54	51.93 \pm 0.33	30.10 \pm 0.69	57.85 \pm 0.68	80.64 \pm 0.43	46.21 \pm 0.89
- sample	34.35 \pm 0.62	52.14 \pm 0.19	27.74 \pm 0.74	56.69 \pm 0.62	79.54 \pm 0.88	45.19 \pm 0.75
gCAS	50.01 \pm 0.62	55.56 \pm 0.59	51.21 \pm 1.74	76.37 \pm 1.60	77.70 \pm 1.46	79.99 \pm 1.03
GDPL	31.89 \pm 0.96	50.14 \pm 0.79	24.99 \pm 1.14	-	-	-
+ sample	34.60 \pm 0.47	45.01 \pm 0.24	31.54 \pm 0.80	-	-	-
DiaAdv	40.97 \pm 0.95	53.44 \pm 0.50	36.84 \pm 1.30	-	-	-
- sample	41.71 \pm 0.47	56.46 \pm 0.45	36.28 \pm 1.48	-	-	-
PEDP	64.63 \pm 0.16	77.03 \pm 1.39	61.77 \pm 1.01	84.12 \pm 0.38	91.66 \pm 0.52	81.19 \pm 0.4
- planning	35.92 \pm 0.54	51.93 \pm 0.33	30.10 \pm 0.69	57.85 \pm 0.68	80.64 \pm 0.43	46.21 \pm 0.89
- ensemble	64.34 \pm 0.29	77.63 \pm 2.04	60.85 \pm 1.54	83.31 \pm 0.55	91.66 \pm 0.78	80.10 \pm 0.55
- sample	66.95 \pm 0.45	78.11 \pm 3.03	65.02 \pm 1.22	84.74 \pm 0.55	92.07 \pm 0.97	81.30 \pm 0.82

Table 2: Standard evaluation results. We report the mean and standard deviation over 5 runs.



Experiments

Dialog pair	Win	Lose	Tie	α
PEDP vs. DiaSeq	41.7	31.3	27.0	0.820
PEDP vs. DiaAdv	36.5	27.6	35.9	0.856
PEDP vs. GDPL	32.6	26.5	40.9	0.839

Table 3: Human evaluation results. We report the mean over 9 judges and Krippendorff's alpha (α) that measures the inter-rater reliability. Typically, results are considered reliable if $\alpha > 0.800$.



Thank you !