# A Pseudo Code

---

**Algorithm 1** CMTA

---

**Initialize**: replay buffer $D$ with $\emptyset$
**Initialize**: initial hidden state $h_0$ with zero tensor
**Initialize**: policy $\pi$ with $\phi$, Q-function Q, task encoder $g$, k experts $f^1, \cdots, f^k$, $lstm$, fully connected layer $\mathcal{W}$
**Input:** state $s_t$ for each environment, one-hot task id $z_\tau$

```
 1: for episode m = 1, 2, · · · do
 2:   for time-step t = 1, 2, · · · do
 3:     for each task τᵢ do
 4:       z_enc^j = f^j(s_t), ∀j ∈ 1, · · · , k
 5:       z_task = g(z_τ)
 6:       h_t = lstm(s_t; h_{t-1})
 7:       α_1, · · · , α_k = softmax(W(h_t; z_task))
 8:       z_enc = Σ_{j=1}^k α_j · z_enc^j
 9:       z = z_task || z_enc
10:       sample action a_t ∼ π(·|z_task; z_enc)
11:       Perform action a_t, get reward r_t and next state s_{t+1}.
12:       D = D ⋃ {s_t, a_t, r_t, s_{t+1}, h_t, z_τ}
13:     end for
14:     randomly sample batch from D
15:     compute L_contrastive by Eq 3
16:     compute L_actor and L_critic by Eq9 and Eq10
17:     update k experts with L_contrastive
18:     update π_φ with L_actor
19:     update all components except π_φ with L_critic
20:   end for
21: end for
```

---

# B Libraries

We use the following open-source libraries: MetaWorld[2], MTEnv[3], MTRL[4].

---

[2] https://github.com/rlworkgroup/metaworld, commit-id:af8417bfc82a3e249b4b02156518d775f29eb289
[3] https://github.com/facebookresearch/ mtenv
[4] https://github.com/facebookresearch/mtrl
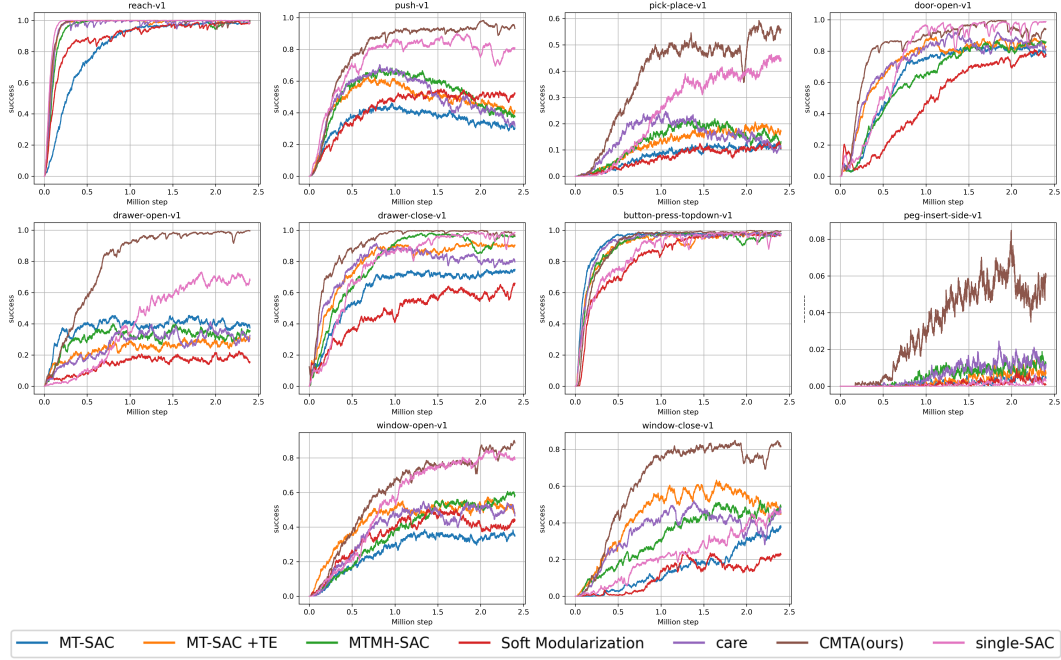
# C   Additional Experiment Results



Figure 6: Training curves of different methods on each task of MT10-Mixed, each curve is averaged over 8 seeds. Our approach consistently outperforms baselines in all tasks, whether on asymptotic performance or sample efficiency.

# D   Hyperparameter Details.

Table 3: Hyperparameter values that are common across all the methods

| Hyperparameter | Hyperparameter values |
| --- | --- |
| batch size | 128 × number of tasks |
| network architecture | feedforward network |
| actor/critic size | three fully connected layers with 512 units |
| non-linearity | ReLU |
| policy initialization | standard Gaussian |
| temperature | learned and distangled with tasks |
| exploration parameters | run a uniform exploration policy 1500 steps |
| num of samples / num of train steps per iteration | 1 env step / 1 training step |
| evaluation frequency | 3000 steps |
| replay buffer size | 5000000 |
| policy learning rate | 3e-4 |
| Q function learning rate | 3e-4 |
| optimizer | Adam |
| policy learning rate | 3e-4 |
| beta for Adam optimizer for policy | (0.9, 0.999) |
| Q function learning rate | 3e-4 |
| beta for Adam optimizer for Q function | (0.9, 0.999) |
| discount | 0.99 |
| Episode length (horizon) | 150 |
| reward scale | 1 |

14

Table 4: Hyperparameter values of task encoder

| Hyperparameter | Hyperparameter values |
|---|---|
| task encoder train from scratch | embedding layer with dim 64 + FC 128 + FC 64 + FC 64 |
| pretrained | pre-trained embedding layer with dim 512 + FC 128 + FC 64 + FC 64 |

Table 5: Hyperparameter values of Soft Modularization

| Hyperparameter | Hyperparameter values |
|---|---|
| task encoder type | train from scratch |
| routing network size | 4 layers and 4 modules per layer with dim 64 |

Table 6: Hyperparameter values of CARE

| Hyperparameter | Hyperparameter values |
|---|---|
| task encoder type | pre-trained embedding layer |
| encoder size | FC 64 + FC 64 |
| number of encoders | 6 |

Table 7: Hyperparameter values of CMTA

| Hyperparameter | Hyperparameter values |
|---|---|
| task encoder type | train from scratch |
| encoder size | FC 64 + FC 64 |
| number of encoders | 6 |
| $\beta$ | 2500 |