

355 **A Proof of Proposition 1**

356 **Proposition 1** (Minimal Representation Insensitive to Policy Bias). *With the Markov chain assumption defined by Eq. (3), for any hidden embedding  $h_k$ , we can derive the upper bound of the  $I(h_k; O)$*   
 357 *where the last term  $I(x; y)$  is a constant with respect to the training process.*  
 358

$$I(h_k; O) \leq I(h_k; x) - I(x; y) \leq I(z; x) - I(x; y) \leq I(z; x), \quad (\text{A.1})$$

359 where the last term  $I(x; y)$  is a constant with respect to the training process.

360 *Proof.* From the Data Processing Inequality (DPI) [18], in this Markov chain, we can obtain

$$I(z; x) \geq I(z; y, O) = I(z; O) + I(z; y|O). \quad (\text{A.2})$$

361 For the second term  $I(z; y|O)$ , suppose  $y$  and  $O$  are independent, we can further factorize it and  
 362 derive

$$I(z; y|O) = H(y|O) - H(y|z, O) \quad (\text{A.3})$$

$$= H(y) - H(y|z, O) \quad (\text{A.4})$$

$$\geq H(y) - H(y|z) \quad (\text{A.5})$$

$$= I(z; y). \quad (\text{A.6})$$

363 As we assume that  $z$  is sufficient, we have  $I(z; y) = I(x; y)$ . Plugging above result back into Eq.(A.2)  
 364 yields

$$I(z; x) \geq I(z; O) + I(z; y|O) \quad (\text{A.7})$$

$$\geq I(z; O) + I(z; y) \quad (\text{A.8})$$

$$= I(z; O) + I(x; y), \quad (\text{A.9})$$

365 which indicates that  $I(z; x) - I(x; y)$  bounds  $I(z; O)$ . And for any hidden embeddings  $h_l$ , according  
 366 to DPI, we have

$$I(h_k; z) \leq I(z; x) \quad \forall k \in \{1, \dots, L\}, \quad (\text{A.10})$$

367 which yields the final result.  $\square$