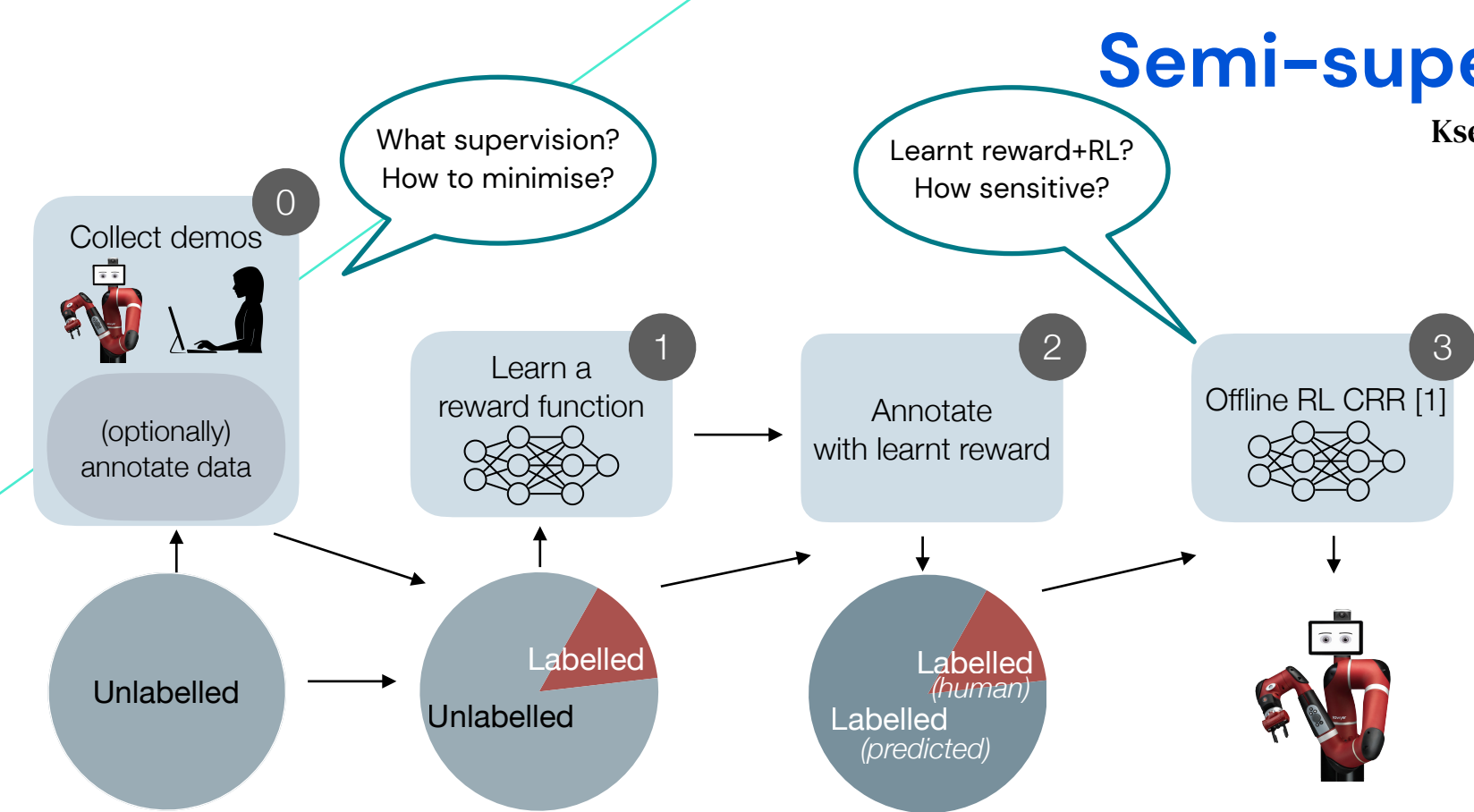


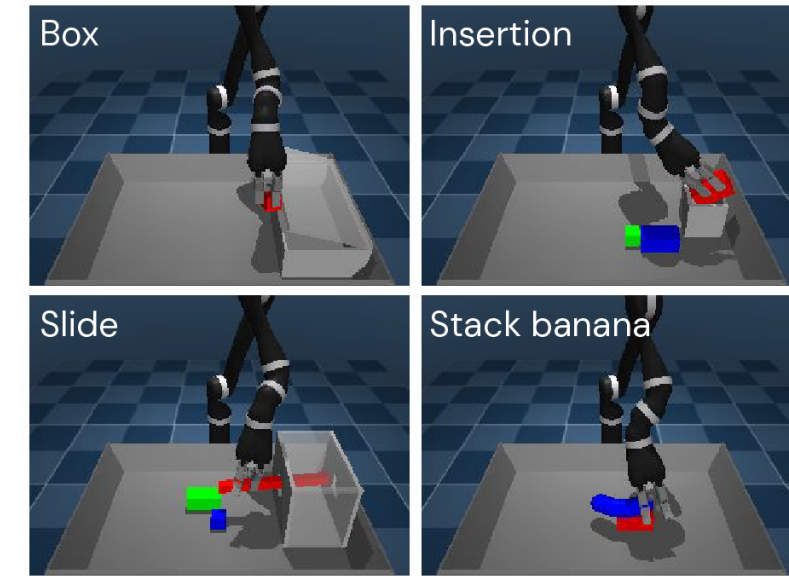
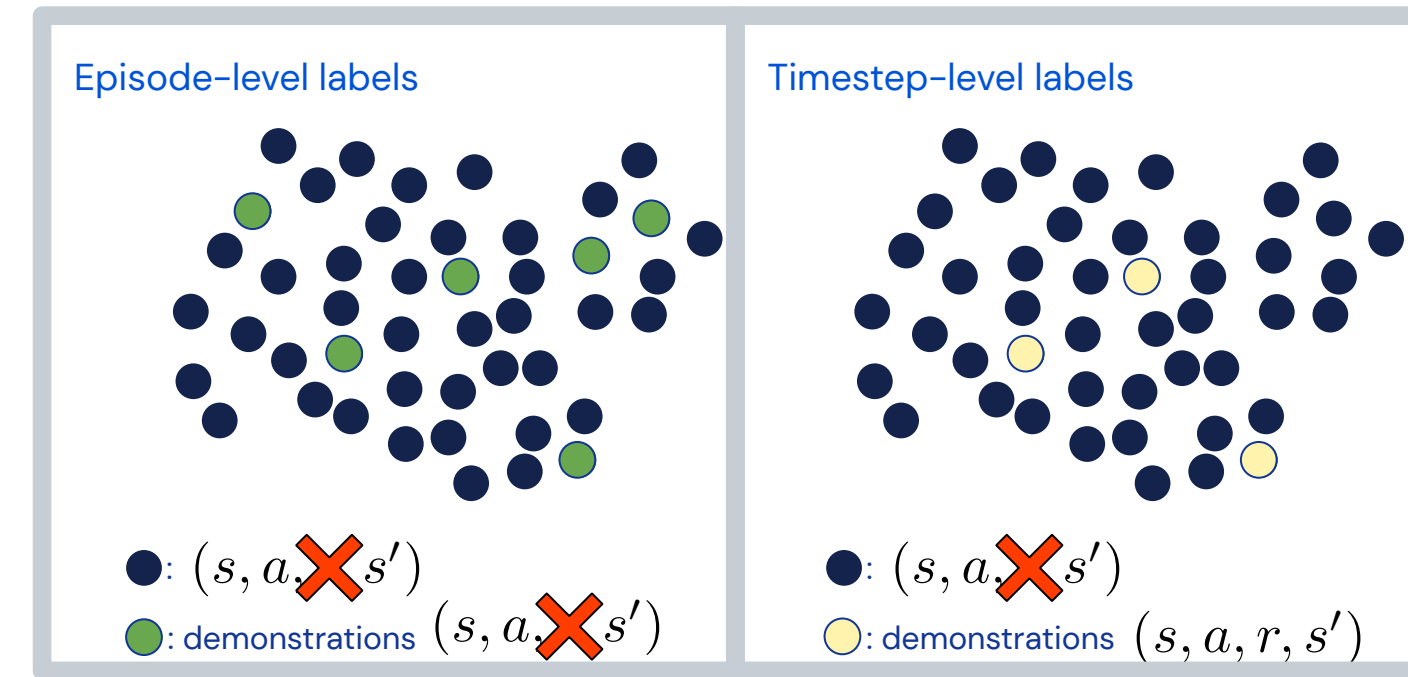
# Semi-supervised reward learning for offline reinforcement learning

Ksenia Konyushova, Konrad Zolna, Yusuf Ayta, Alexander Novikov, Scott Reed, Serkan Cabi, Nando de Freitas

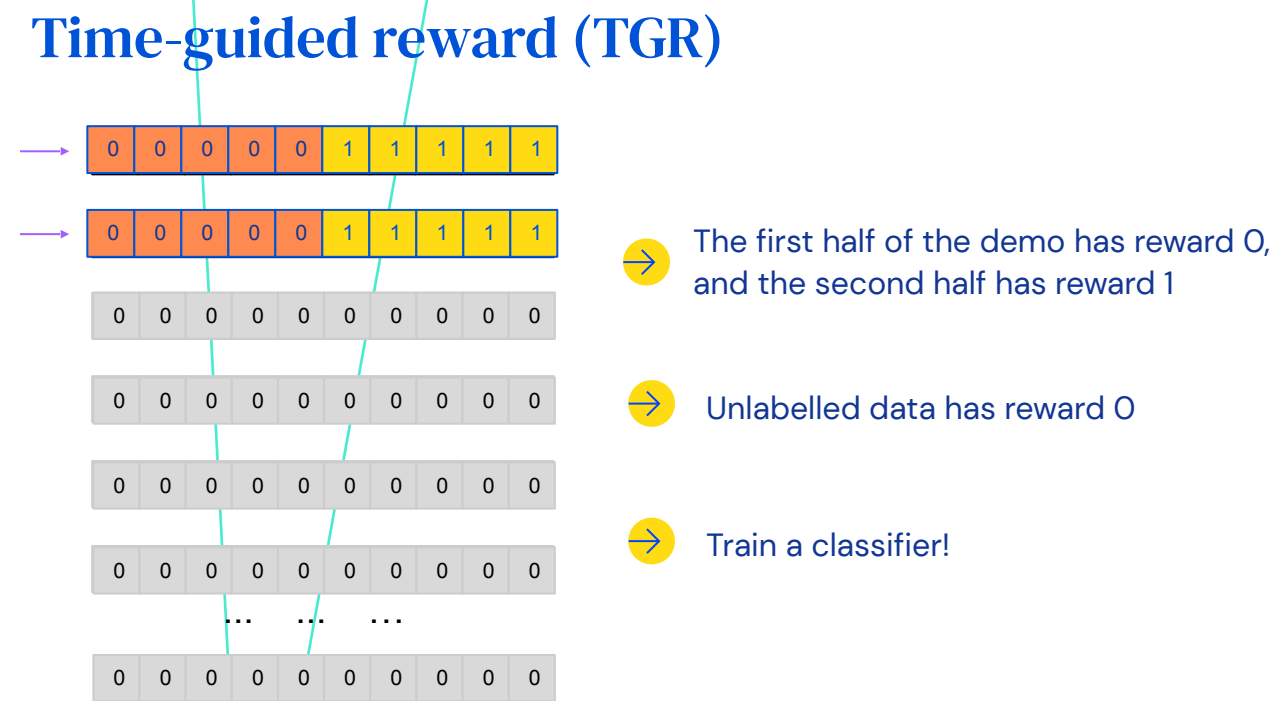
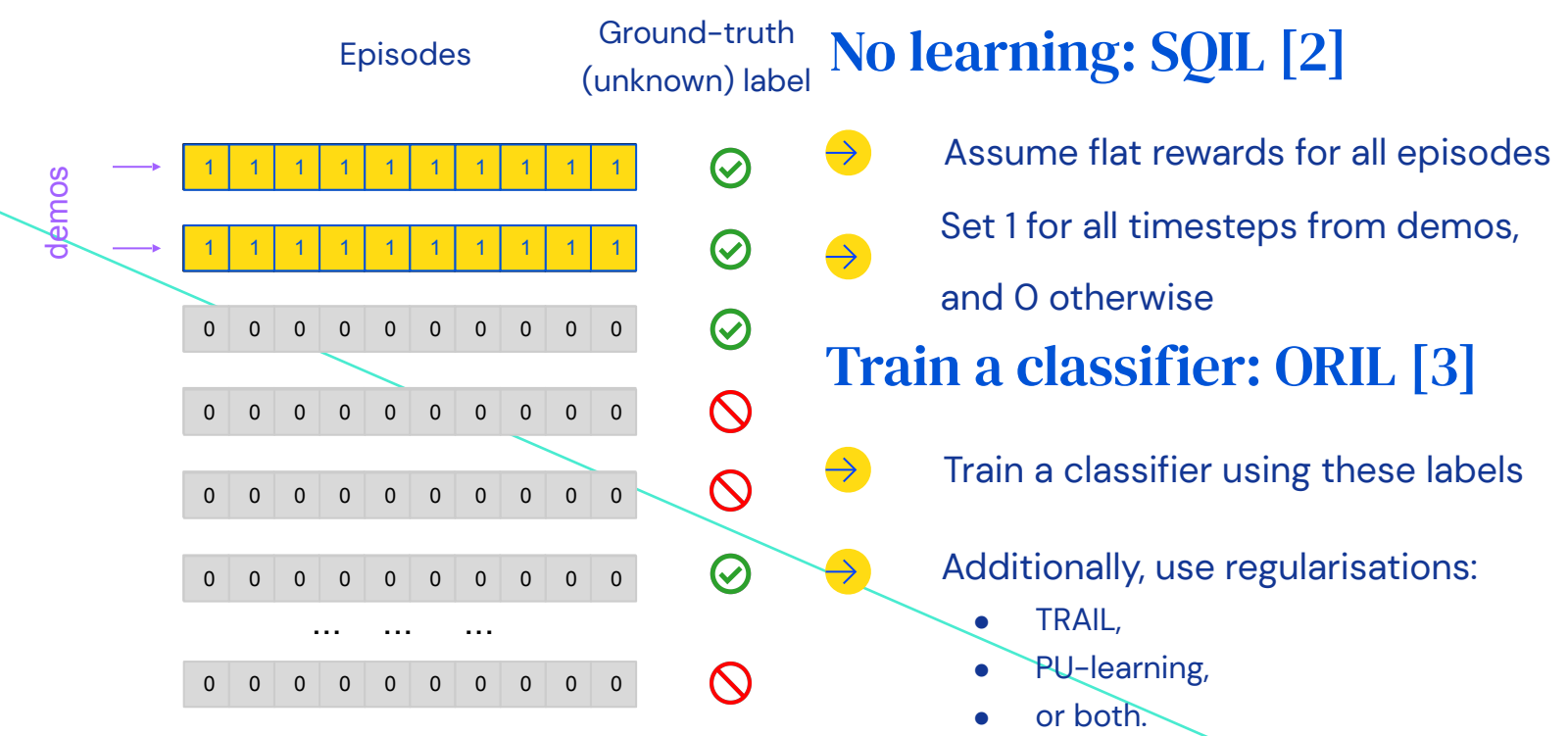


	Problem	Solution
Data	<ul style="list-style-type: none"> <li>hard to collect</li> <li>hard to explore</li> <li>data-hungry</li> </ul>	<ul style="list-style-type: none"> <li>use offline RL</li> </ul>
Reward	<ul style="list-style-type: none"> <li>no reward in practice</li> <li>hard to engineer</li> <li>impossible for some tasks</li> </ul>	<ul style="list-style-type: none"> <li>learn reward function</li> </ul>

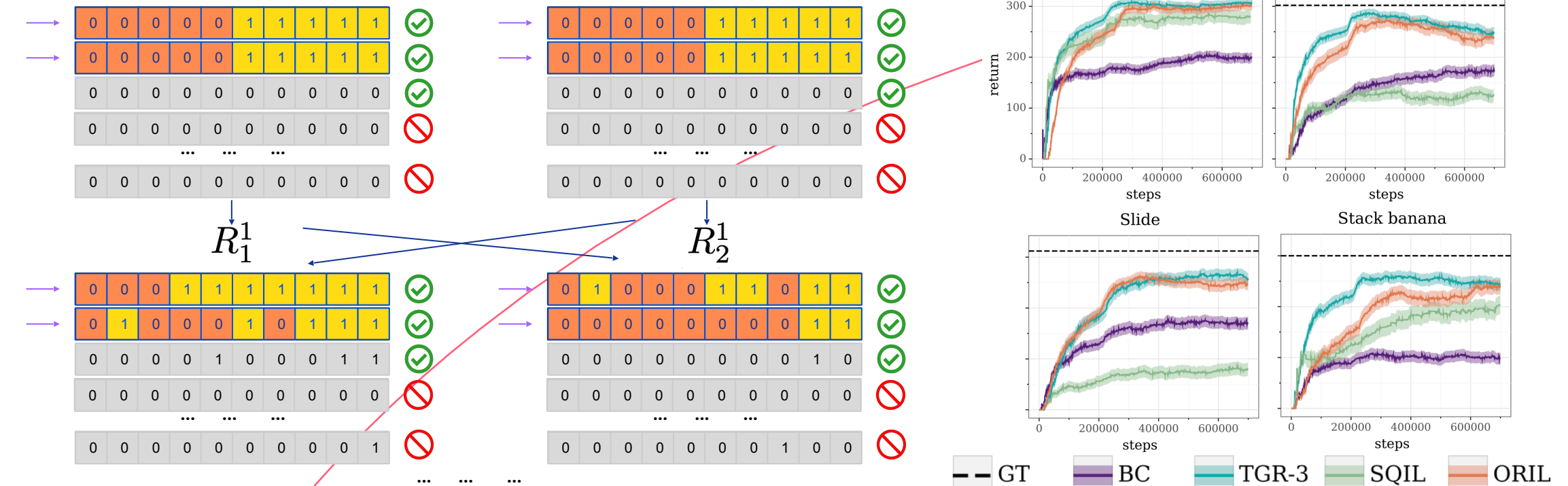
Needs human annotations, try to minimise



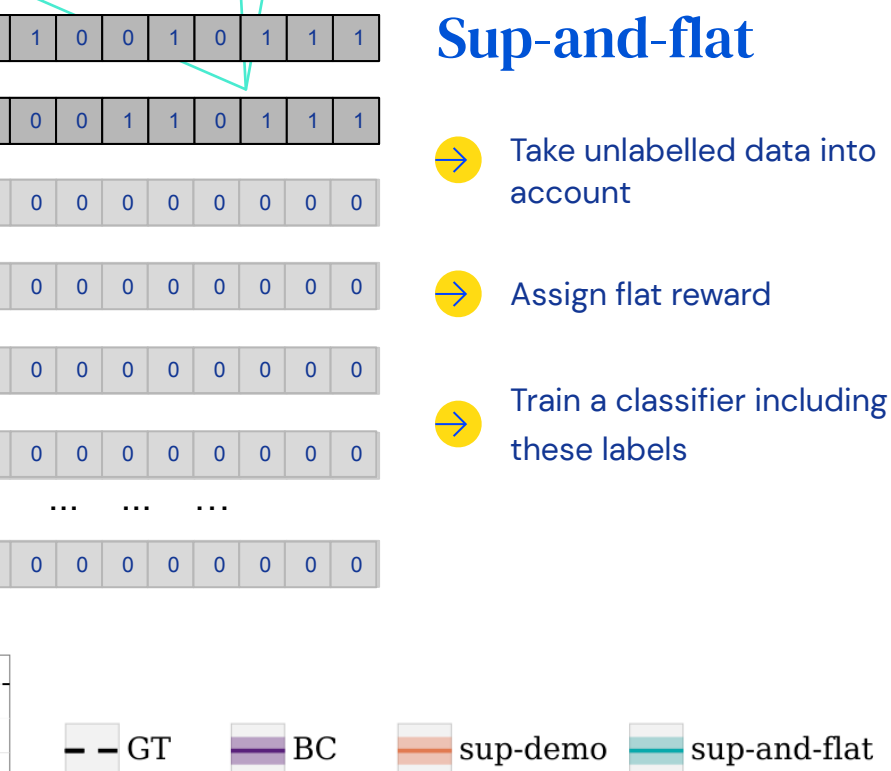
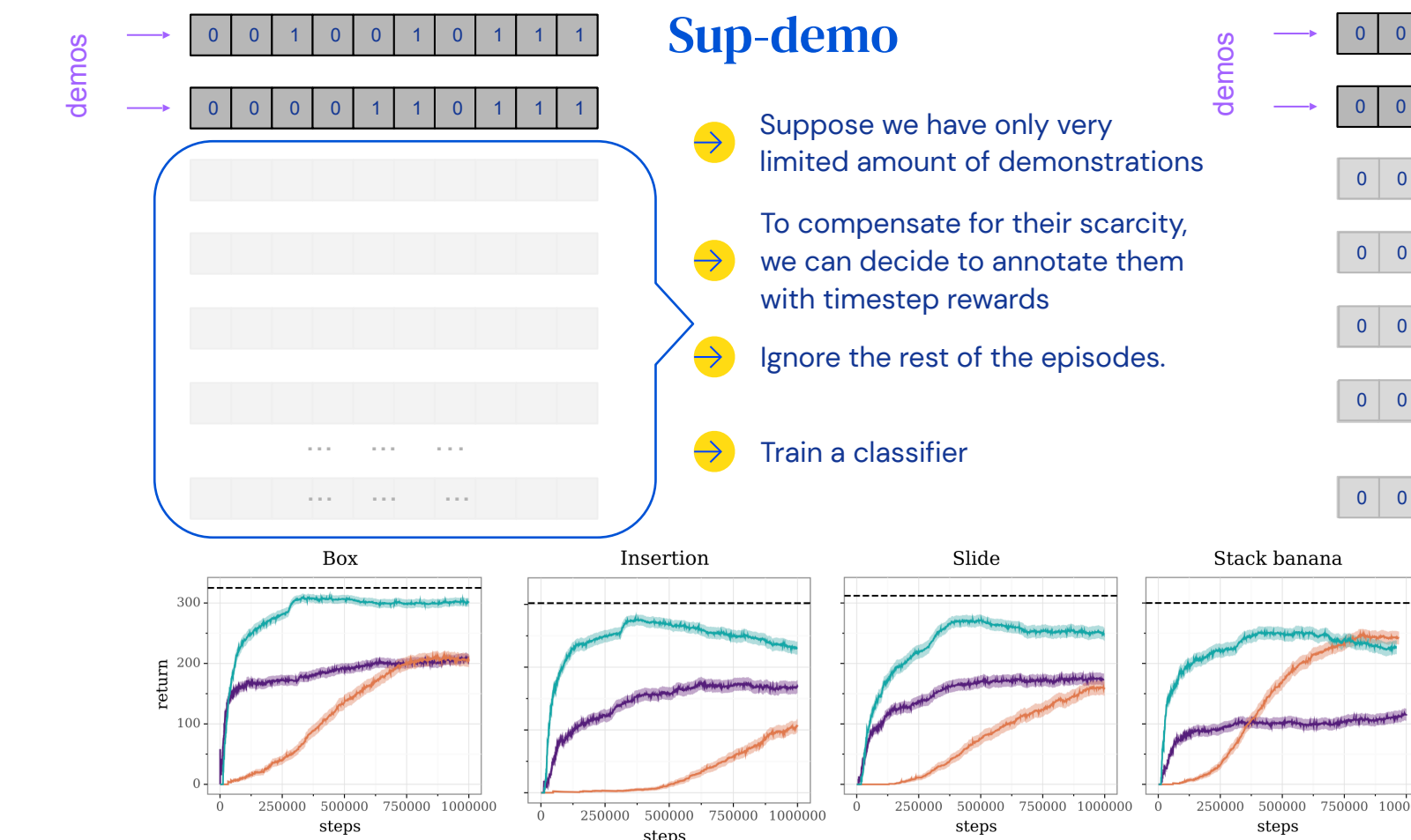
## Episode-level labels



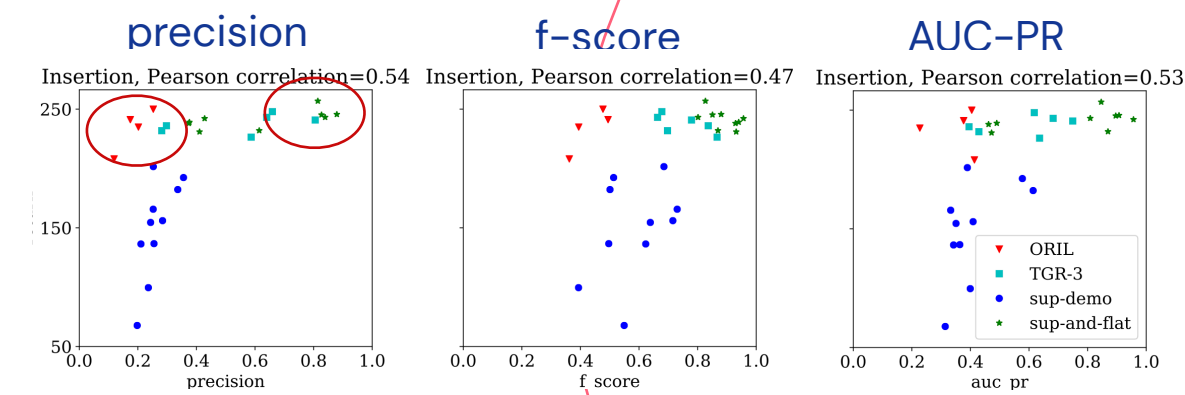
## TGR improvement



## Timestep-level labels



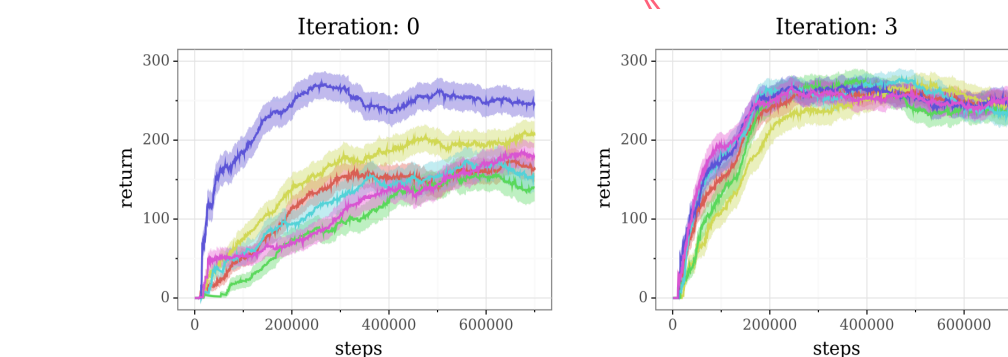
## How to choose the best reward model?



## Other RL algorithms?



## TGR improvement reduces the variance



[1] Critic Regularized Regression. Z. Wang, A. Novikov, K. Zolna, J.T. Springenberg, S. Reed, B. Shahriari, N. Siegel, J. Merel, C. Gulcehre, N. Heess, N. de Freitas, NeurIPS 2020

[2] SQIL: imitation learning via reinforcement learning with sparse rewards. S. Reddy, A. D. Dragan, and S. Levine. In ICLR, 2020.

[3] Offline Learning from Demonstrations and Unlabeled Experience. K. Zolna, A. Novikov, K. Konyushkova, C. Gulcehre, Z. Wang, Y. Ayta, M. Denil, N. de Freitas, S. Reed.