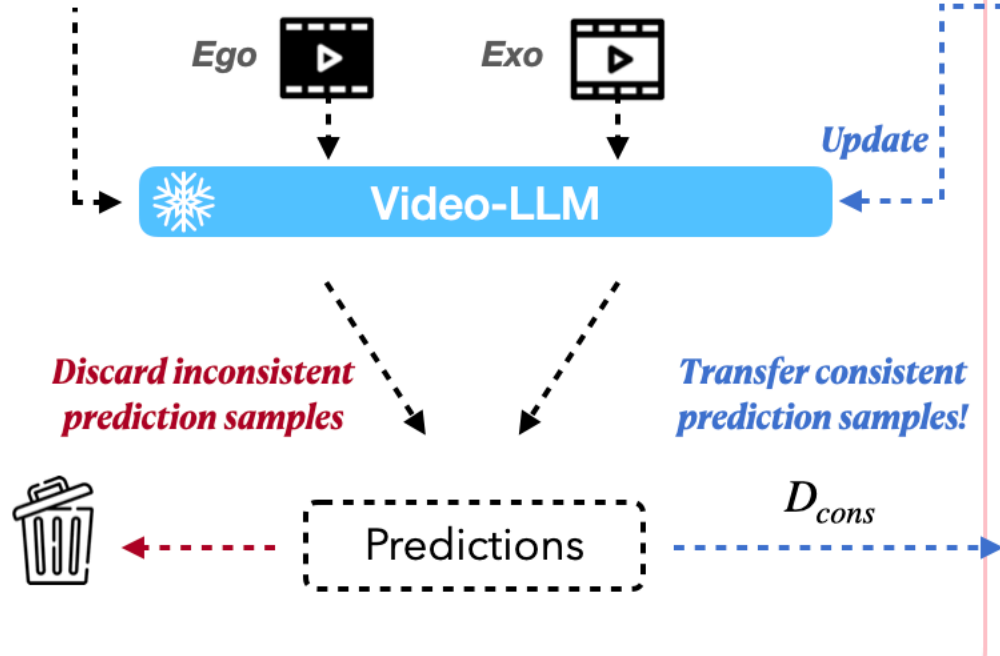


Self-Guided Consistent Sample Mining

Training Set D

Q: Does <event> happen from <timestamps> in the video?



Dual-Encoder architecture

