| Stage 1 | Stage 2 & 3 |
|---|---|
| Attribute-wise | All-attribute |

SFT

SFT

Model

Inference: Selected Attribute $\mathcal{A}'$

Video + $\mathcal{A}'$ ➡ ASID-Captioner ➡ Caption