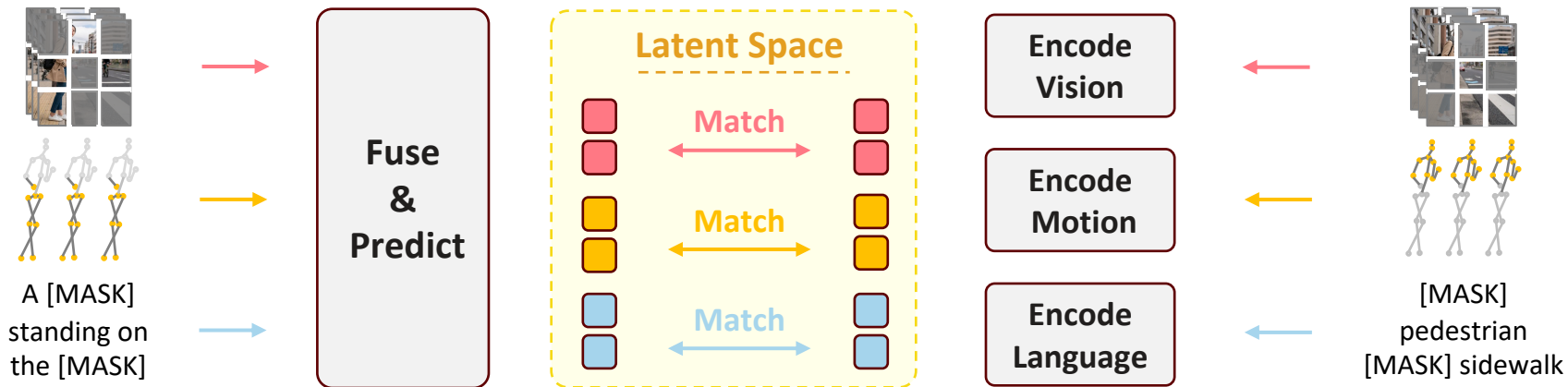


Pretraining Stage



Downstream Tasks

Action Recognition

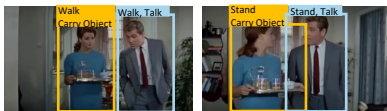
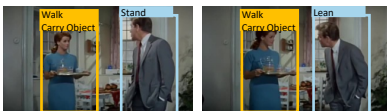


Action:
Running



Action:
Play Tennis

Action Localization

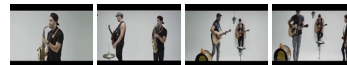


Text-Video Retrieval



Caption: a man throws an American football at an aiming board

Video QA



Question: how many guys perform a song together?

Answer: Three