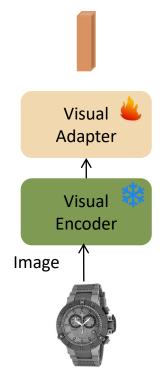
Stage 1: Feature Alignment

Text paired with image

A grey watch with an army style strap.

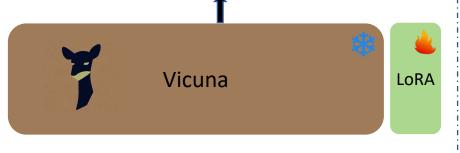




Stage 2: Boundary Perception

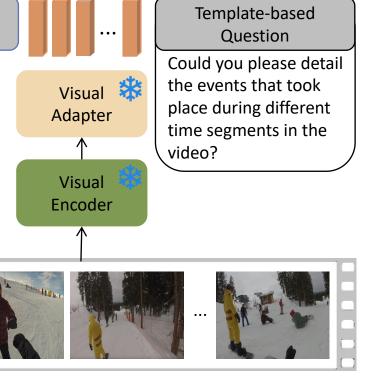
Template-based Answer

People are snowboarding down a large hill of snow, from 00 to 93. The people get to the bottom and start taking their snowboards off, from 94 to 99.





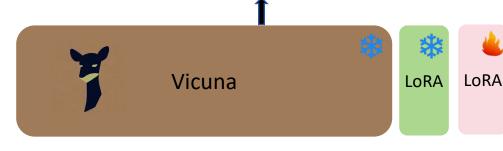
Video



Stage 3: Instruction Tuning

High-quality Answer

Of course! These two people play rock-paper-scissors twice. The first time is from 50 to 65, and after the woman puts a blindfold on, they play again, from 81 to 88.

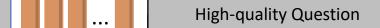


Visual 🗱

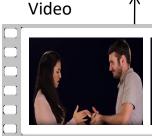
Adapter

Visual

Encoder



Can you help me count how many times the two people play rock-paper-scissors in the video?



System

Prompt



