

Fusion of Captioning

"video_caption": "A person is holding a blender with their hand",

}

"task_procedures": "<person> is <verb> holding <action> a blender with their hand."