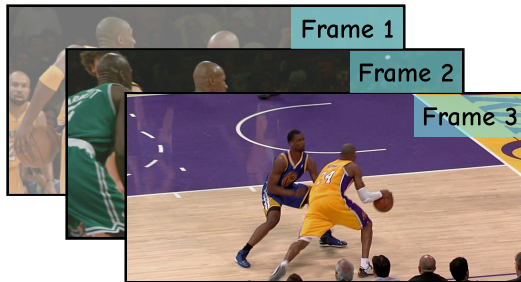


## Uniform Segmentation Frame



## Captions of each Frame

Caption 1

Caption 4

Caption 2

Caption 5

...

Caption 3

Caption n

Video Caption

Static Role-playing

Finetune



Question-  
Answers Pairs



Inference

Summary  
Context