

T consecutive frames from a single shot

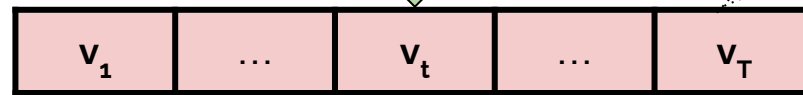
Frame 0

Frame t

Frame T

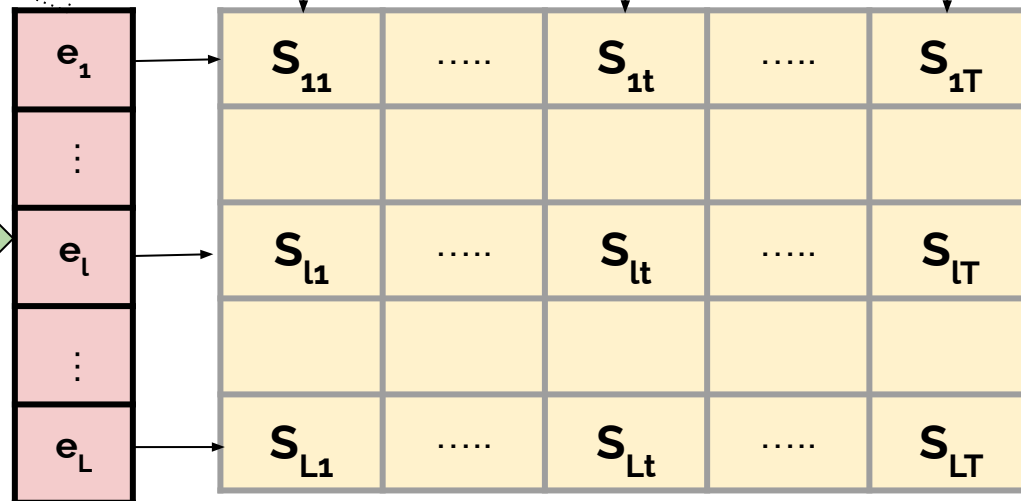


CLIP
visual
encoder



512 D frame wise embeddings

512 D label wise embeddings



Temporal Pooling

Temporal Pooling

Temporal Pooling



CLIPSceneScore₁

Threshold

Living room,
Lounge

label

Living room
...
Lounge
...
Stadium

Curated labels from
scene taxonomy

A photo of a
{label}, a type
of
background
location

CLIP text
encoder