

(۱۱) توضیح سیم

در این شبکه به نظر می رسد که یک سیم video captioning است. و به این صورت عمل می کند که ابتدا یک لایه 3D conv روی تصاویر می زد و ویژگی ها و متناوب با دوری می بصری را استخراج می کند و این عملیات را برای هر می فریم سیم می انجام می دهد. سپس فریم لایه 3D conv را به یک CNN می دهد و در نهایت یک لایه mean pooling این ویژگی ها را ترکیب می کند. در واقع در این جا به کمک 3D conv ویژگی های که به طول زمان ارتباط دارند را استخراج می کند و پس برادر ویژگی استخراج شده را به عنوان ورودی CNN می دهد و در نهایت چون فریم های سیم به هم اطلاعات می برساند و قابل ترکیب دارند که یک لایه mean pooling می توانی اطلاعات را استخراج می کند.