

Multimodal Abstractive Summarization for Open-Domain Videos

Jindřich Libovický¹, Shruti Palaskar², Spandana Gella³, Florian Metze²



¹Charles University

²Carnegie Mellon University

³University of Edinburgh



CHARLES UNIVERSITY



THE UNIVERSITY
of EDINBURGH

Introducing Summarization with How2 Data

- Summarization
 - Present subset of information in a more compact form (maybe across modalities)
- “Description” field
 - 2-3 sentences of meta data: template based, uploader provides
 - “Informative” and abstractive summary of a how-to video
 - Should generate interest of a potential viewer



How To Make a Spanish Omelet : Cutting Peppers for A Spanish Omelet

1,307 views



Published on Mar 4, 2008

How to cut peppers to make a Spanish Omelette; get expert tips and advice on making traditional Cuban breakfast recipes in this free cooking video.

SUBSCRIBE 3.3M

How2 Dataset

- 2000h of **how-to** videos
 - 300h for MT
 - 480h for ASR
 - 80K videos
 - Object, Scene and Action features
- Human annotated transcripts
- Meta-data
- Video summaries / descriptions
- Very different topics
 - Cooking, music, sports, etc.



How To Make a Spanish Omelet : Cutting Peppers for A Spanish Omelet

1,307 views

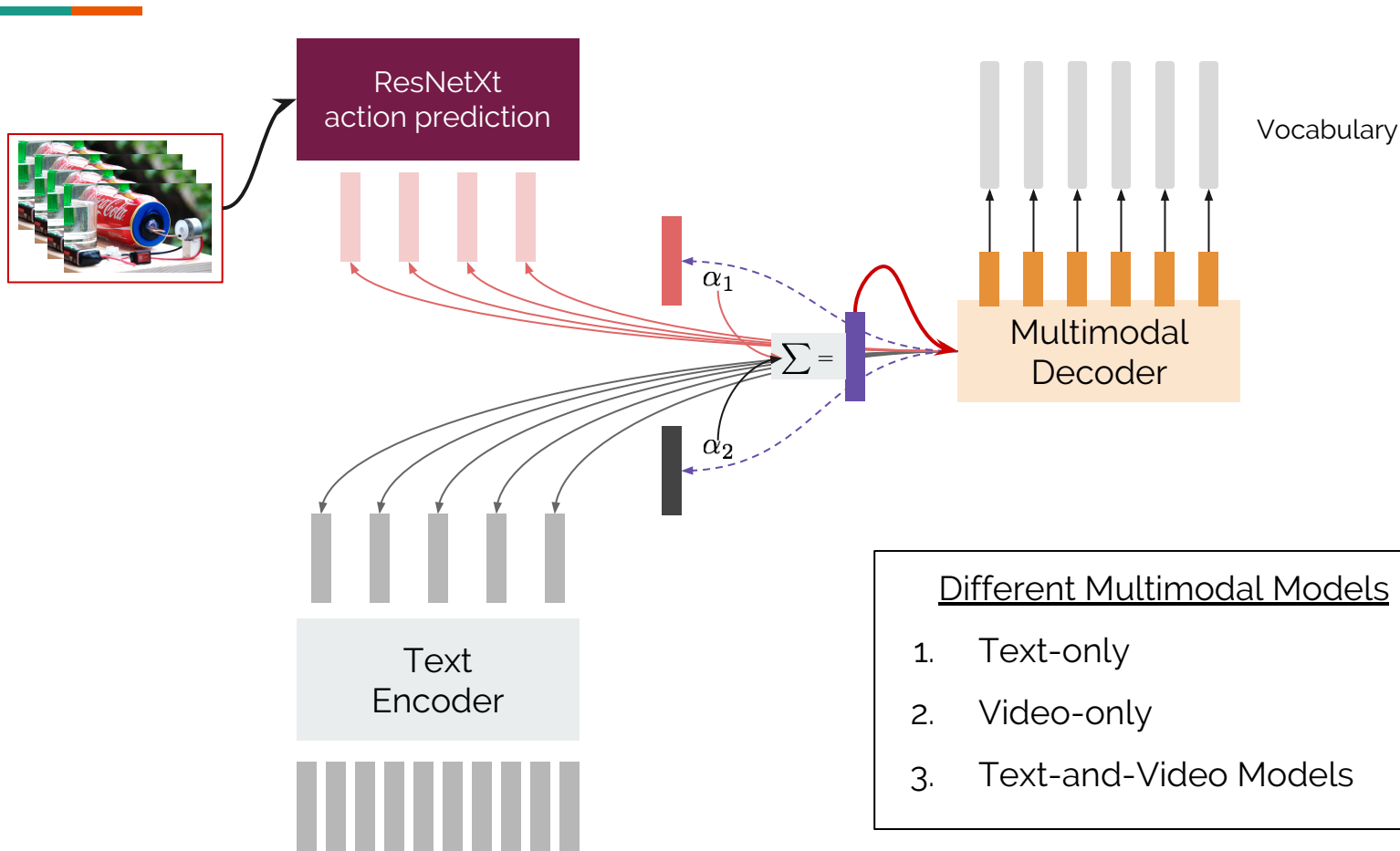


SUBSCRIBE 3.3M

How to cut peppers to make a Spanish Omelette; get expert tips and advice on making traditional Cuban breakfast recipes in this free cooking video.

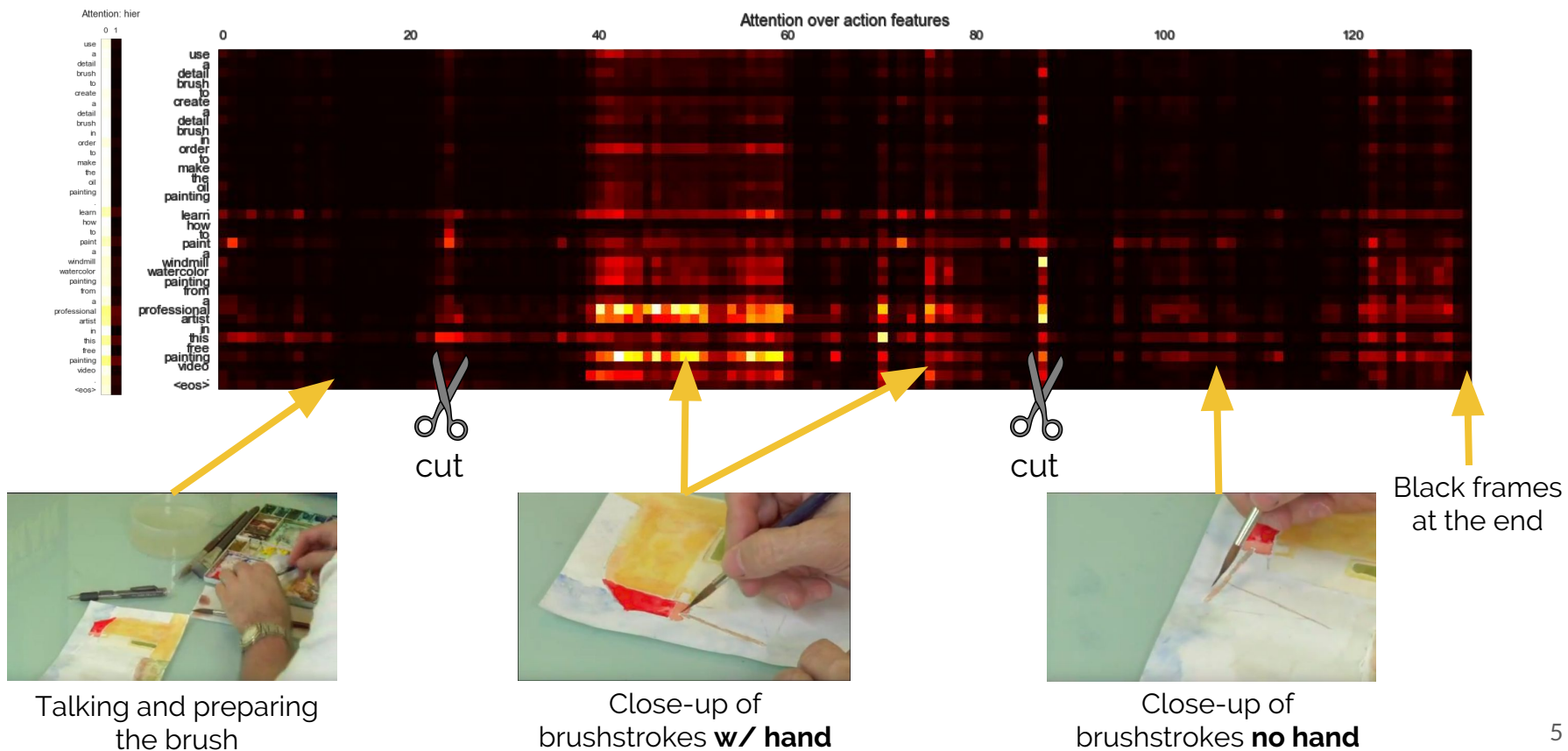
R. Sanabria et al., "How2: A Large-Scale Dataset for Multimodal Language Understanding"
ViGIL 2018

Hierarchical Multi-modal Attention



- Different Multimodal Models
1. Text-only
 2. Video-only
 3. Text-and-Video Models

Attention over the Video Features



Winning Model in DSTC7 AVSD Track



- Applied this model to the DSTC7 AVSD track
- Goal was to test the effect of pre-training for this task on a large corpus, How2
- Charades dataset: Visual Question Answering type data with audio and summary as extra modalities
- Best performing model in objective and human evaluation

Visit our posters!



- ❖ Multimodal Abstractive Summarization for Open-Domain Videos
- ❖ How2: A Large-Scale Dataset for Multimodal Language Understanding

Dataset: <https://github.com/srvk/how2-dataset>

Code: <https://github.com/lium-lst/nmtpytorch>

How2 Dataset

nmtpytorch