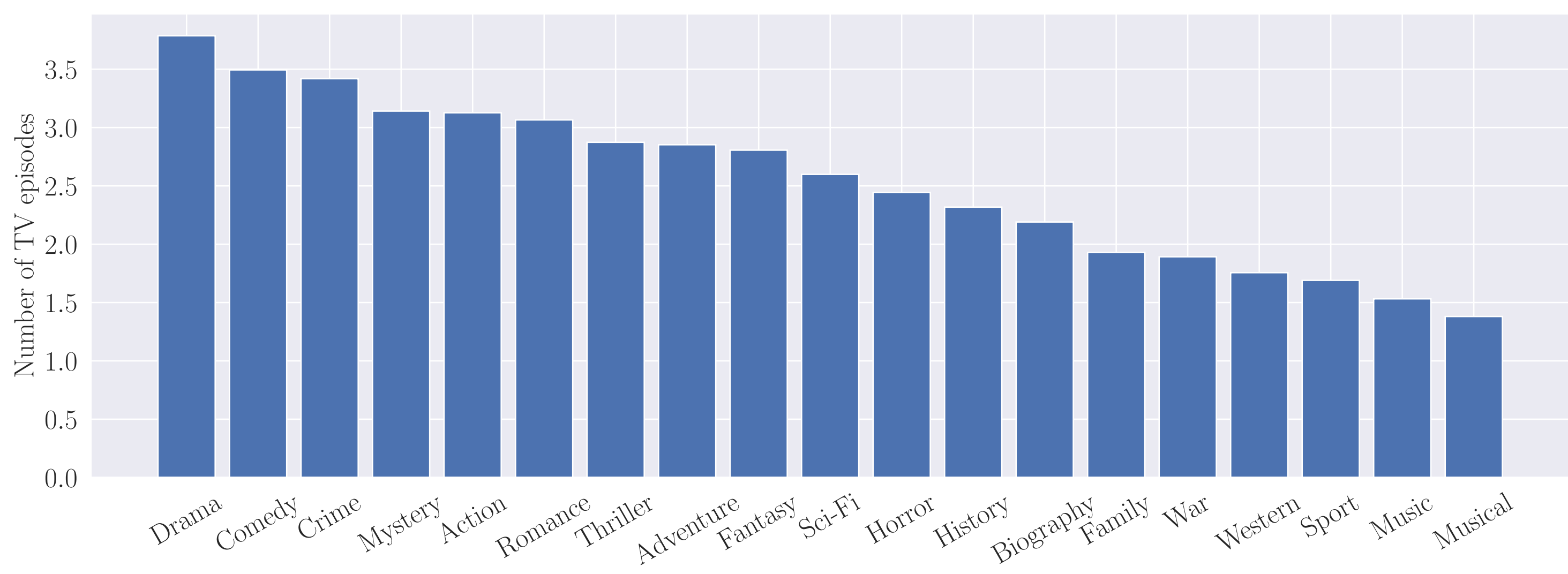


ONE Large-scale VSTAR Dataset

Dataset	Vision	Language	Scene	Topic	# Dialogues	# Turns	Turns/Clip	Words/Turn
VisualDialog	Image	QA	\times	\times	120K	2.4M	20.0	4.0
Twitch-FIFA	Live Video	Dialogue	\times	\times	15K	161K	10.4	6.0
AVSD	Recorded Video	Dialogue	\times	\times	11K	118K	20.0	9.5
MoiveNet [†]	Movies	Dialogue	\checkmark	\times	-	421K	-	7.2
OpenViDial 2.0	Movies & TV Series	Dialogue	\times	\times	116K	5.6M	48.0	8.3
VSTAR (Ours)	TV Series	Dialogue	\checkmark	\checkmark	185K	4.6M	25.1	6.7



	Scene	Video	Source
OVSD	300	21	MiniFilm
BBC	670	11	Documentary
MovieScenes	21k	150	Movie
MovieNet	42k	318	Movie
VSTAR (Ours)	265k	8159	TV Episode

	Sentence	Sent/Seg	Language
DiaSeg_711	19K	5.6	English
Doc2Dial	19K	3.5	English
ZYS	12K	6.4	Chinese
VSTAR (Ours)	4.6M	9.3	English

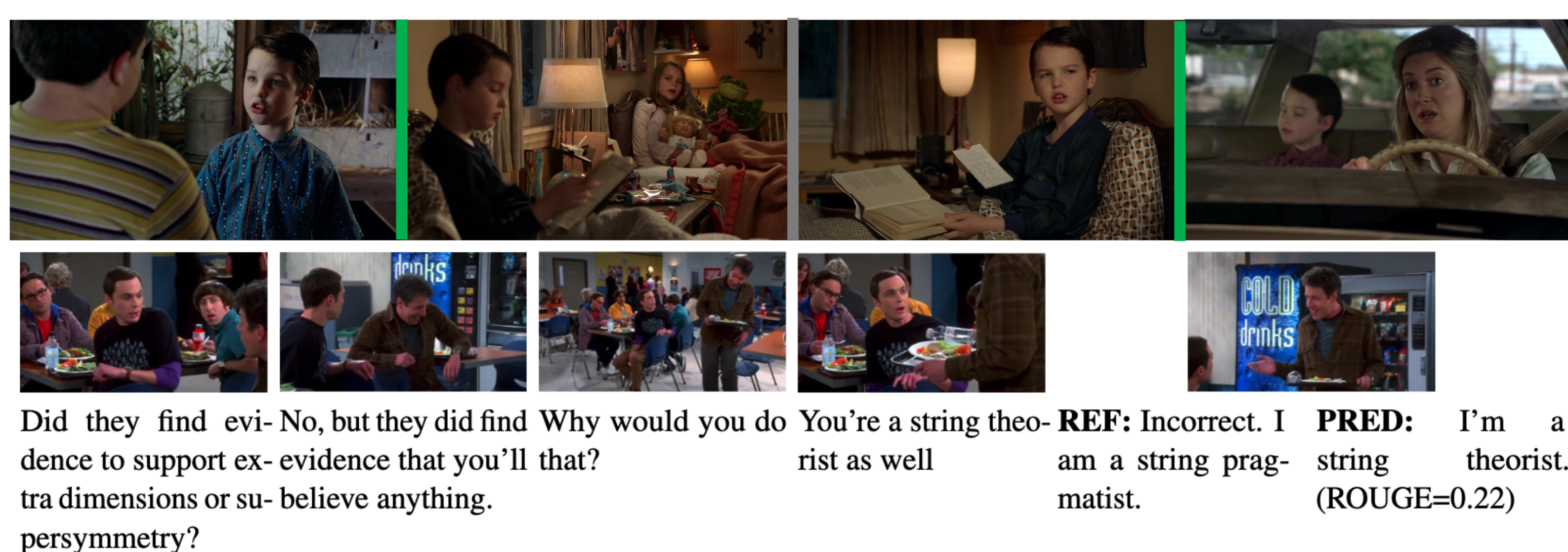
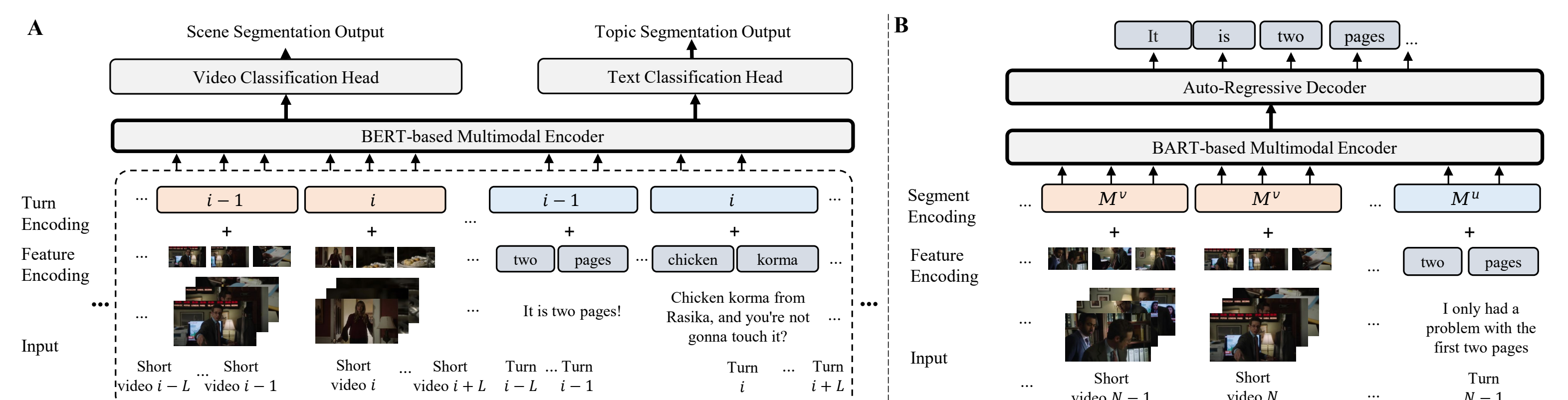
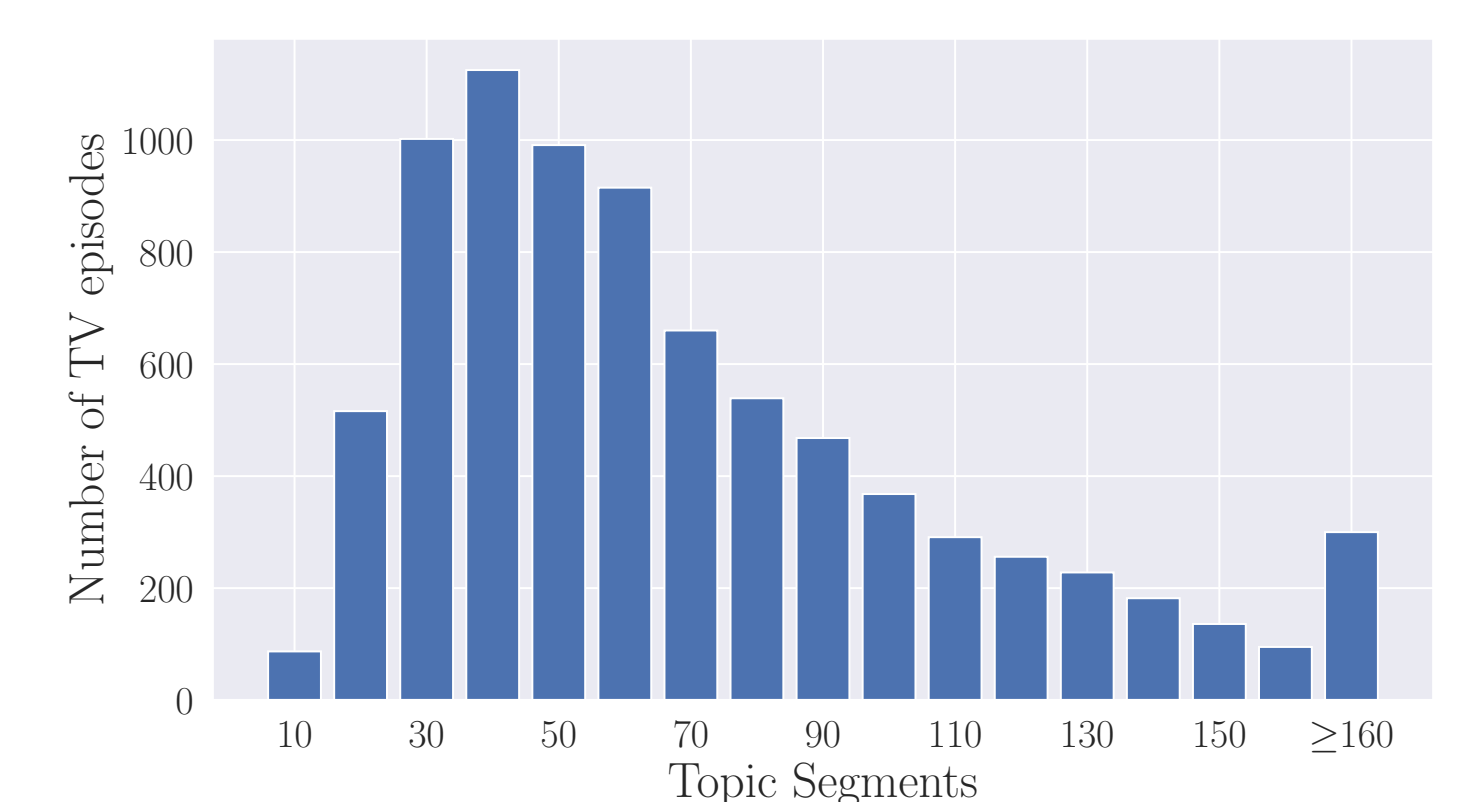
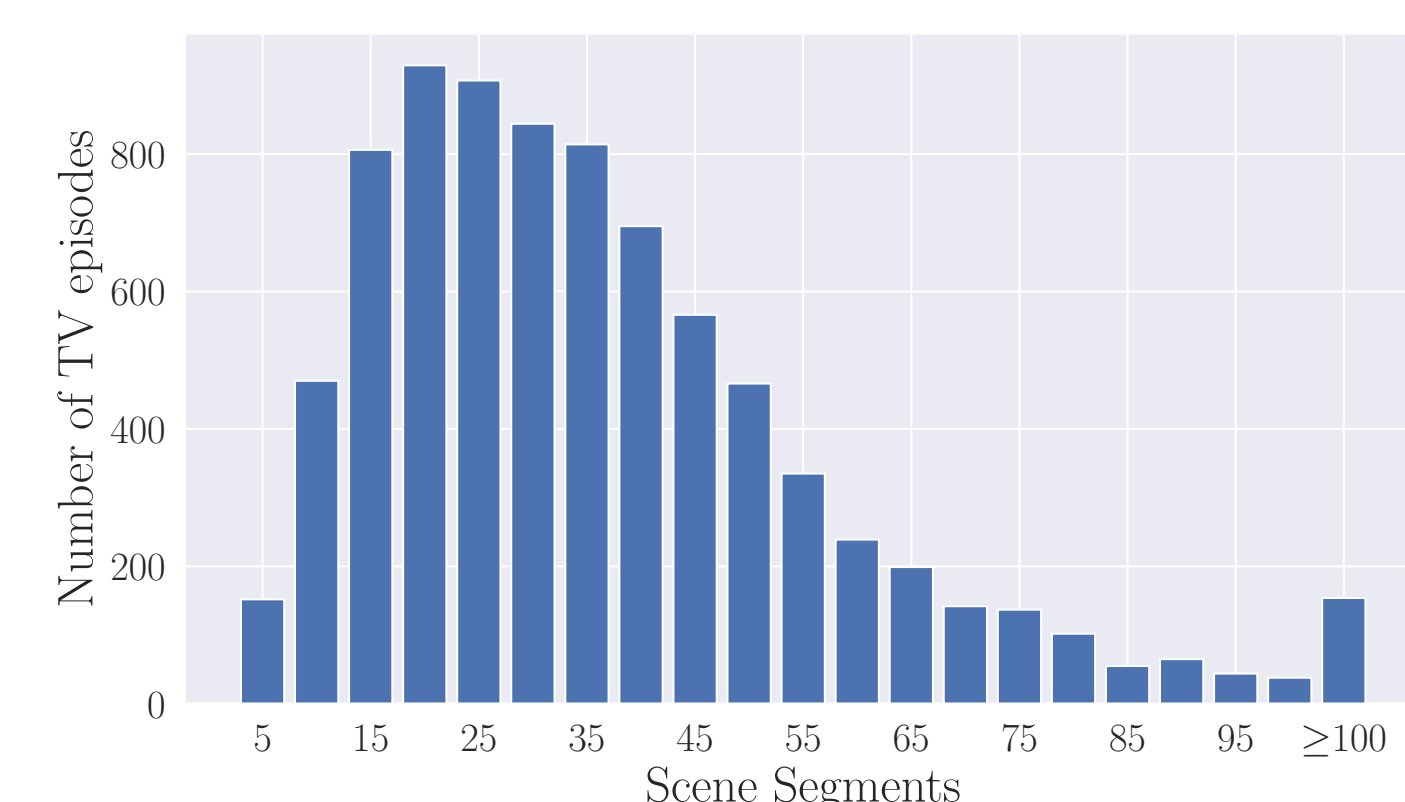


Figure 4. Cases for segmentation and generation tasks. The green bars indicate scene transition, and the grey bars indicate topic transition. **PRED** denotes our predicted response. **REF** denotes the reference human response.