

On the Challenges in Evaluating Visually Grounded Stories

Aditya K Surikuchi Raquel Fernández Sandro Pezzelle



UNIVERSITY OF AMSTERDAM

Institute for Logic, Language and Computation

Visually Grounded Story Generation

Visually Grounded Story Generation

E#ž 3 _ S` VgU]eXadLahVž VdMSdWTg^VfZaVē[fZW S^` VsdZ[t S` VZW
LahVēZ[eXSUM[XVSč: WfQVēfa YVfSi Skl Tgf S` afZVd bVēa` YdSTeZ[_
Xca_ TVZ[Vž: WZaVeZ[_ [b^SUMSeFZW_ S` efgYYVēfa YVfSi Skž
EgWWW^kl fZWtg^VfēfSd^ k[Yž 4afZ _ W VgU] adV[hWēdeSXVtkž 8aaV S` V
T[feaXb^SeVdYa k[YSeFZWtg^Vfē ^fZWēba_ žFi a_ W SdNZ[fi S` V fZVē
TafZ Ya k[YTSU]ž

Visually Grounded Story Generation

E#ž 3 _ S` VgUj eXadLahVz` VdMSdWTg^VfZaVd[fZW S^` VsdZ[t S` VZW
LahVdZ[eXSUN[XScž: WfVdVefa YVfSi Skl Tgf S` afZVdbVda` YdSTeZ[_
Xda_ TVZ[Vž: WZaVeZ[[b^SUNSeFZW_ S` efdgYYVdVefa YVfSi Skž
EgVWV'kt fZWtg^VfVefSd[k[Yž 4afZ _ W VgUj adV[hVadeS^Vž 8aaV S` V
T[feaXb^SeVdYa k[YSeFZWtg^VfV ^fZWda_ ž Fi a_ W SdNZ[fi S` V fZVd
TafZ Ya k[Y TSU] ž

E\$ž;` fZWabg^Wf daa_ t S_ S` [S fgj Wa i Seb^Sk[Y fZVdUbdV[a t
aT[h[agefa fZVdZSaeg` XaV[Y Scdg` VZ[ž EgVWV'kt S` afZVd_ S` Tgdef
[fa fZWdaa_ t Z[eXSUNda` fadW i [fZ S` YVdSeZWg` YW Sf fZW def_ S` ž W
eVd` V_ S` eS SU] i Se_ Vfi [fZ Sei [Lbg` fVdXda_ fZW def i Za_ S` SYW
fa bgeZ Z[TSU] ž 3efZVd efdgYYVf S UagV aXVgef ^W fZVd[d aTelgd[Y
fZVd_ ahW Wfež;` fZW [Vef aXfZWda_ af[a t S fZ[dv_ S` SbbVsdVf Z[e
VbdVd[a` a` VdXezau] S` V XsdSeZW SflZW fZWfi a_ W YdSbbVf

Visually Grounded Story Generation

3 _ S` VgUj eXadLahVz VdMSdMTg^MfZaVef fZW S^ VsdZL t S` VZW
LahVdz[eXSUM^ XScz: WfQVefafa YVfSi Skl Tgf S` afZVdbVda` YdSTeZL
Xda_ TVZ[Vž: WZaVeZL [b^SUMSefZW S` efcgYYVefafa YVfSi Skž
EgVWV'kt fZWTg^MfVefSdf k[Yž 4afZ _ W VgUj adV[hWadeSXVikž 8aaV S` V
T[feaXb^SeVdYa k[YSe fZWTg^MfVef ^fZWcha_ ž Fi a_ W SdNZ[fi S` V fZVW
TafZ Ya k[Y TSU] ž

;` fZWabg^Wf dda_ t S_ S` [S fgj Wa i Seb^Sk[Y fZWSUhdV[a t
aT[h[agefa fZWUZSaeg` XaV[Y Scbg` VZL ž EgVWV'kt S` afZVd_ S` Tgdef
[fa fZWcha_ t Z[eXSUMUa` fadW i [fZ S` YVdSeZWg` YW Sf fZW def_ S` ž W
eVd` V_ S` eS SU] i Se_ Vfi [fZ Sei [Ubg` fVdXda_ fZW def i Za_ S` SYW
fa bgeZ ZL TSU] ž 3efZVW efcgYYVf S UagV aXVgef ^W fZWS[d aTelgd[Y
fZV[d_ ahW Wfež;` fZW [Vef aXfZWUa_ af[a` t S fZ[dv_ S` SbbVsdM Z[e
VbdVef[a` a` VdXezau] S` V XsdSeZW SflZW fZWfi a_ W YdSbbVf

Datasets

H; EF

Sequences constructed using images from Flickr albums.

Lacks consistency of entities (e.g., *US*, *SU*, *AF*, *W*, *aT*, *W*, *e*)

Corresponding stories are generally descriptive in nature

HI B

Sequences constructed using scenes from movies

Semantically well-connected with recurring characters

Stories contain diverse entities, are longer, and coherent

Models



Qwen-VL



LLaVA



DeepSeek-VL

General-purpose VLMs
(a ~~ZZVZVZ~~)

TAPM (+LLAMA 2)

Specific to Visual Story
Generation

Evaluation

Human evaluation is challenging in terms of:

- Scalability and costs

- Selecting qualified annotators and reliability

Evaluation

Human evaluation is challenging in terms of:

- Scalability and costs

- Selecting qualified annotators and reliability

Automatic Metrics:

- BLEU, METEOR, CIDEr, SPICE, ROUGE

Evaluation

Human evaluation is challenging in terms of:

- Scalability and costs

- Selecting qualified annotators and reliability

Automatic Metrics:

~~BLEU, METEOR, CIDEr, SPICE, ROUGE~~

Reference-free metrics that assess stories along different aspects—

Evaluation

Human evaluation is challenging in terms of:

- Scalability and costs

- Selecting qualified annotators and reliability

Automatic Metrics:

~~BLEU, METEOR, CIDEr, SPICE, ROUGE~~

Reference-free metrics that assess stories along different aspects—

d. γ measure that computes distances between individual metric scores of model- and corresponding human stories

Results

What is overall best performing model?

Results

What is overall best performing model?

How did the models fare along each of the 3 dimensions?

"SFUIFSF PUIFS EJNFOTJPOT SFMFW
WJTVBMMZ HSPVOEFE TUPSJFT

