# An animated picture says at least a thousand words: Selecting Gif-based Replies in Multimodal Dialog

Xingyao Wang (xingyaow@umich.edu) & David Jurgens (jurgens@umich.edu)

UNIVERSITY OF MICHIGAN

BLAB LAB

## 1. What are we trying to do?

Online conversations include more than just text, people like using reaction gifs in their messages.

However, current NLP dialog systems (e.g., chatbots) are almost all text-based.

To fill the gap, we introduce a **new task:** Select a gif-based reply to a text message from a user.

**@user:** Ahhhh! The deadline is in 24 hours!

↳ **@model:**



## 2. New dataset for multimodal text-gif dialog!

1.56M text-gif conversation turns & Metadata for 115k Gifs.

Metadata includes annotated tags (~⅓ of gifs, e.g. happy, exciting), extracted captions. and machine-predicted object names and features (e.g., face, building)

**Object names:** face woman

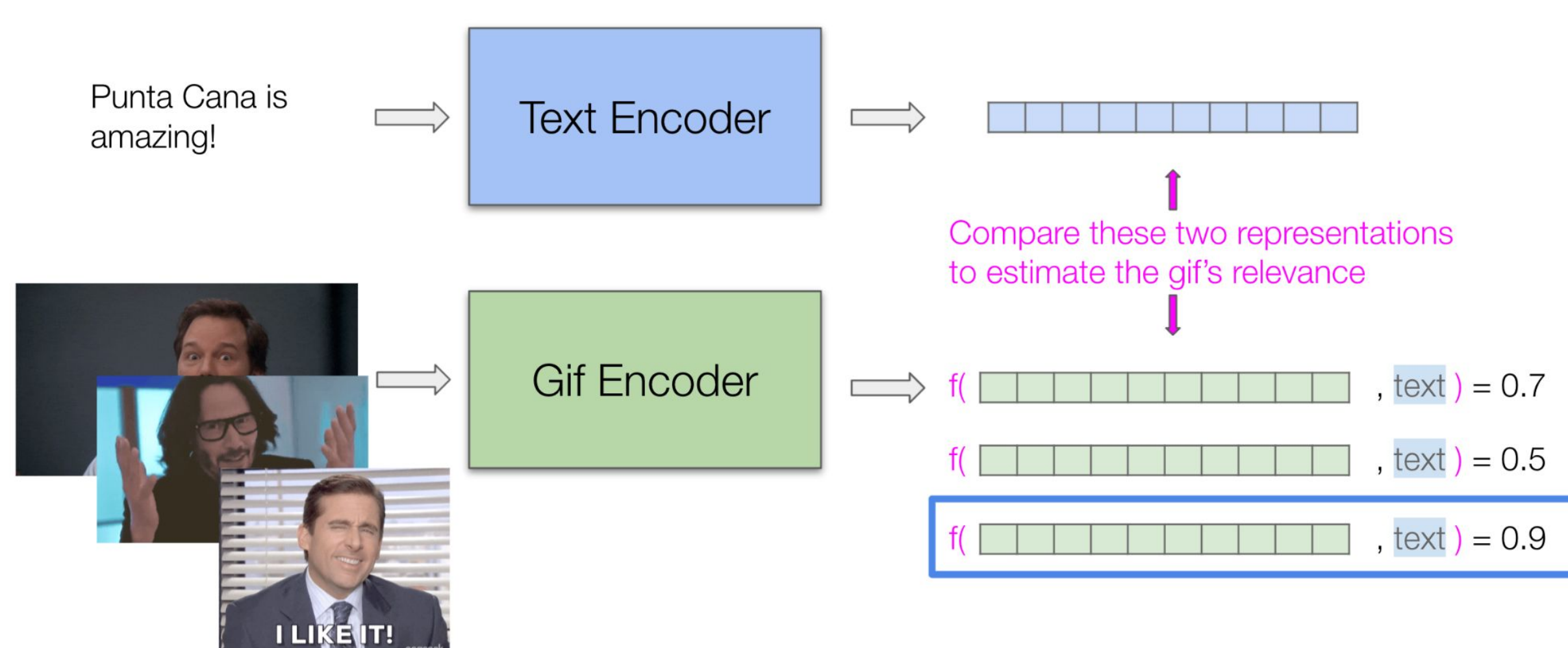**Tweet:** @USER is my hero

↳**Reply:**
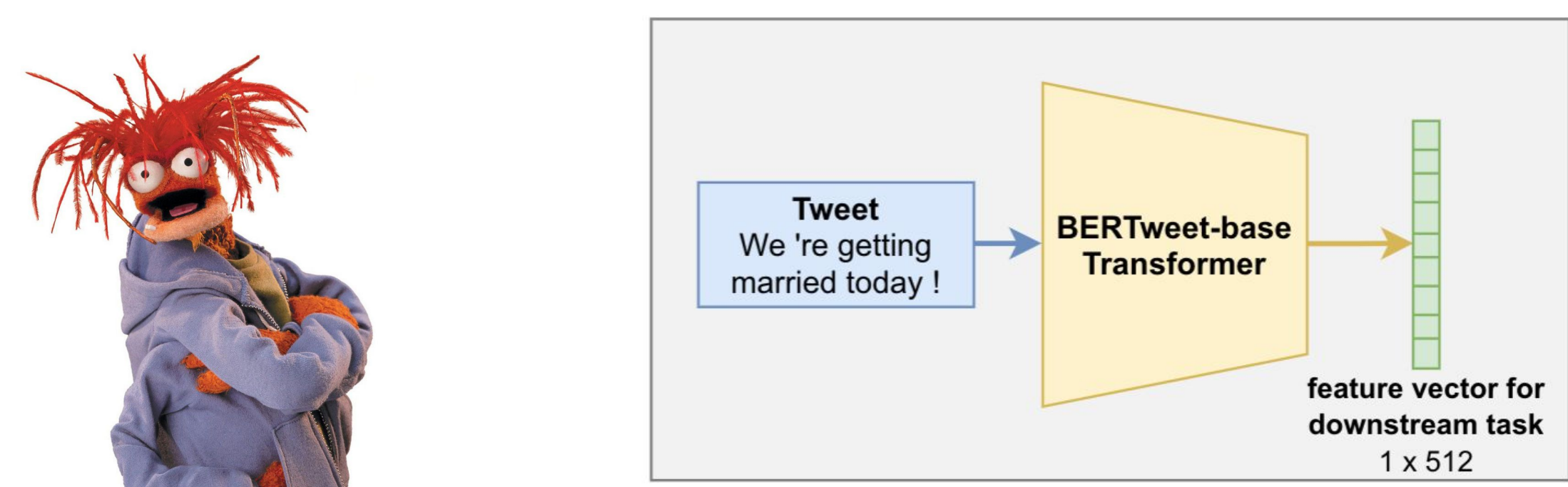


↳ **Annotated Tags:** ["thank"]

↳ **Captions:** Aww , thank you

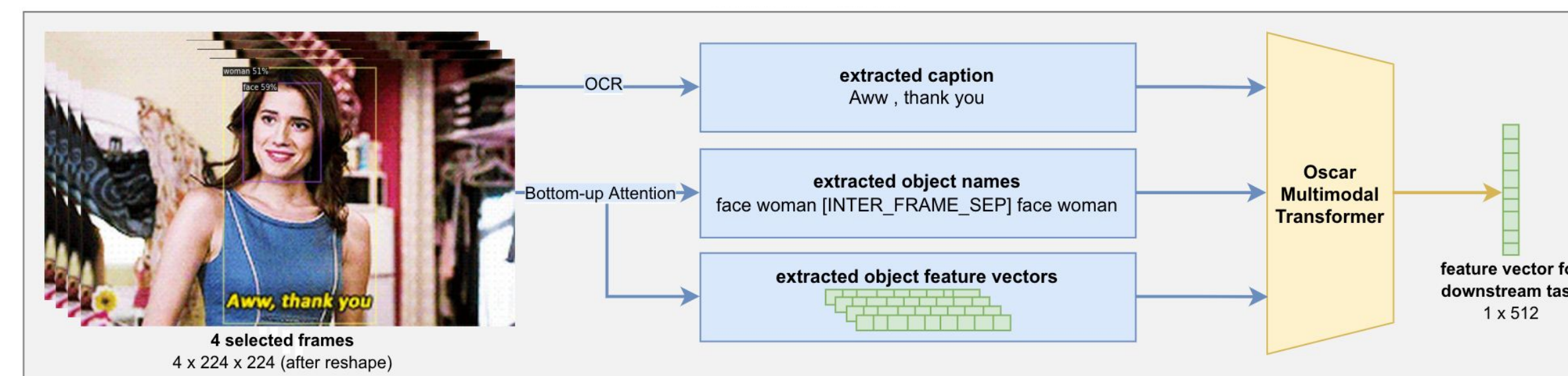## 3. How to select a gif reply? Ranking approach!

We compare **text** and **gif** representations to estimate the gif's relevance. We pick the gif with the highest score as reply.



Punta Cana is amazing! → Text Encoder →

Compare these two representations to estimate the gif's relevance

Gif Encoder →
f( , text ) = 0.7
f( , text ) = 0.5
f( , text ) = 0.9

## 4. Our Model: Pepe the King Prawn[1]



**Tweet** We 're getting married today ! → BERTweet-base Transformer → feature vector for downstream task 1 x 512

Text encoder: *Roberta* (pre-trained on Twitter)



4 selected frames 4 x 224 x 224 (after reshape)

OCR → extracted caption Aww , thank you

Bottom-up Attention → extracted object names face woman [INTER_FRAME_SEP] face woman

extracted object feature vectors

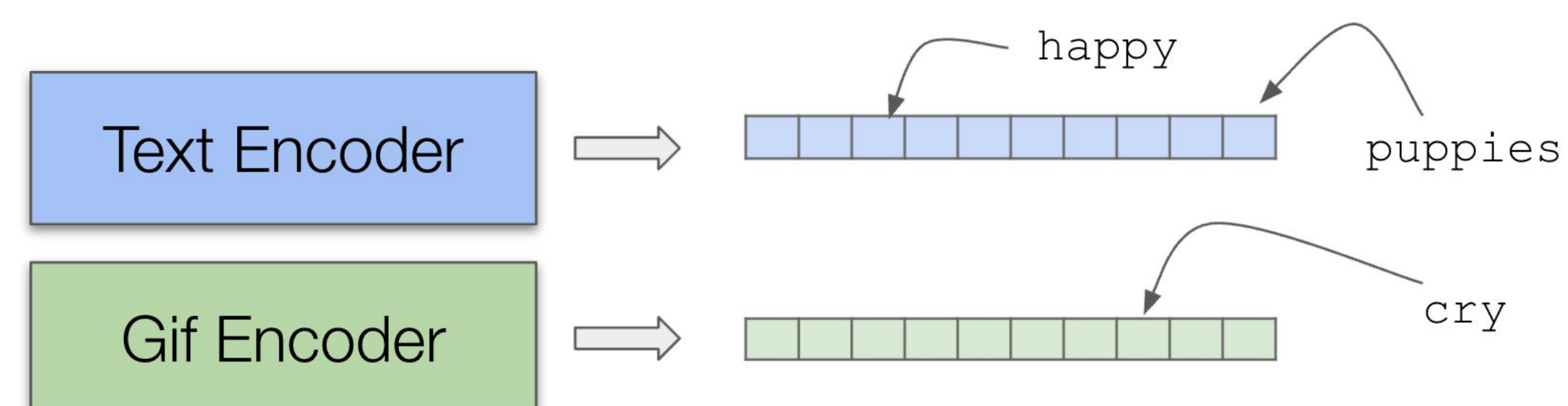→ Oscar Multimodal Transformer → feature vector for downstream task 1 x 512

GIF Encoder: *Oscar Multimodal Transformer*

GIF encoder fuses information about the gif from text modality (captions), visual modality (object features), and bridge between the two modality (object names).
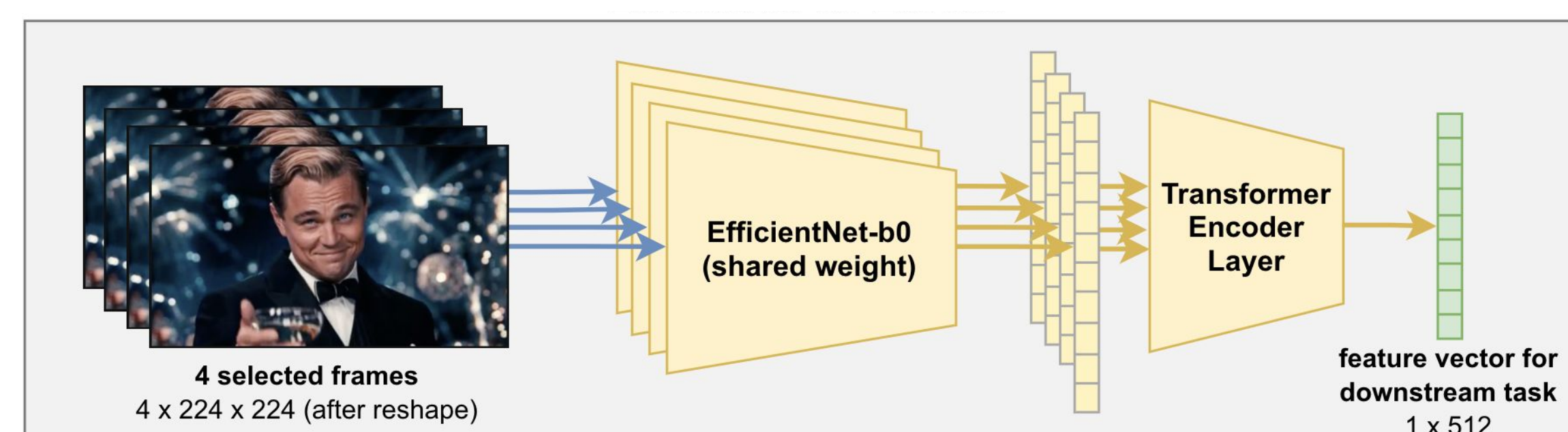
[1] "King Prawn" refers to "selecKting INteresting Gifs for Personal RespAWNses"

## 5. Can we rank gifs differently?

**Idea 1 (Tag-based):** Some gifs have describe content/intent, so use a **tag-based encoder** for both gif and text



Text Encoder → happy puppies

Gif Encoder → cry

**Idea 2 (CLIP-based):** Simplify the gif encoder to use a CLIP-like **CNN-based Image encoder**.



4 selected frames 4 x 224 x 224 (after reshape) → EfficientNet-b0 (shared weight) → Transformer Encoder Layer → feature vector for downstream task 1 x 512
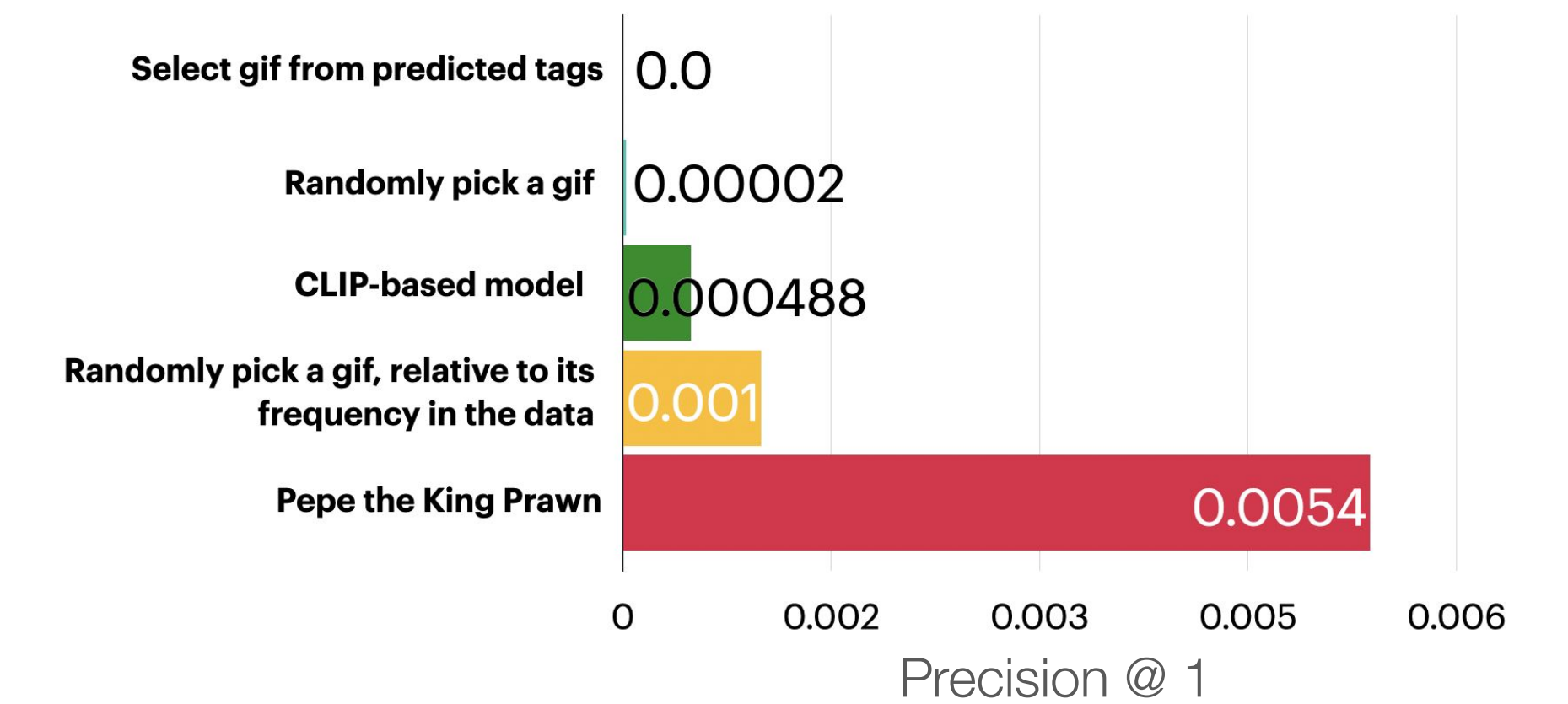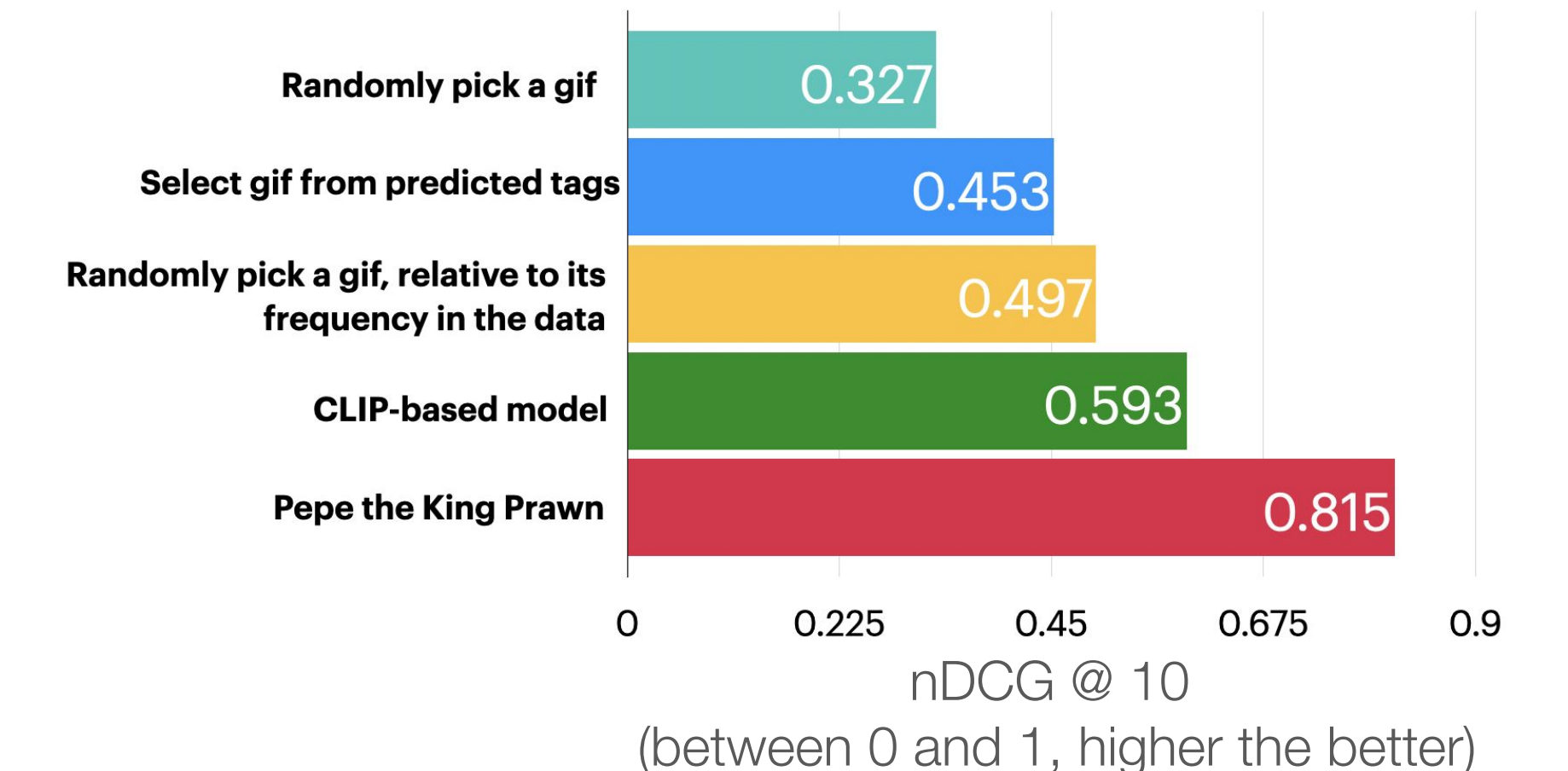
**Other baseline models:**
- **Random model**: randomly select a Gif ID
- **Distribution sampling**: sample Gif ID with probability proportional to its frequency (i.e. more frequently-used gif are sampled more often)
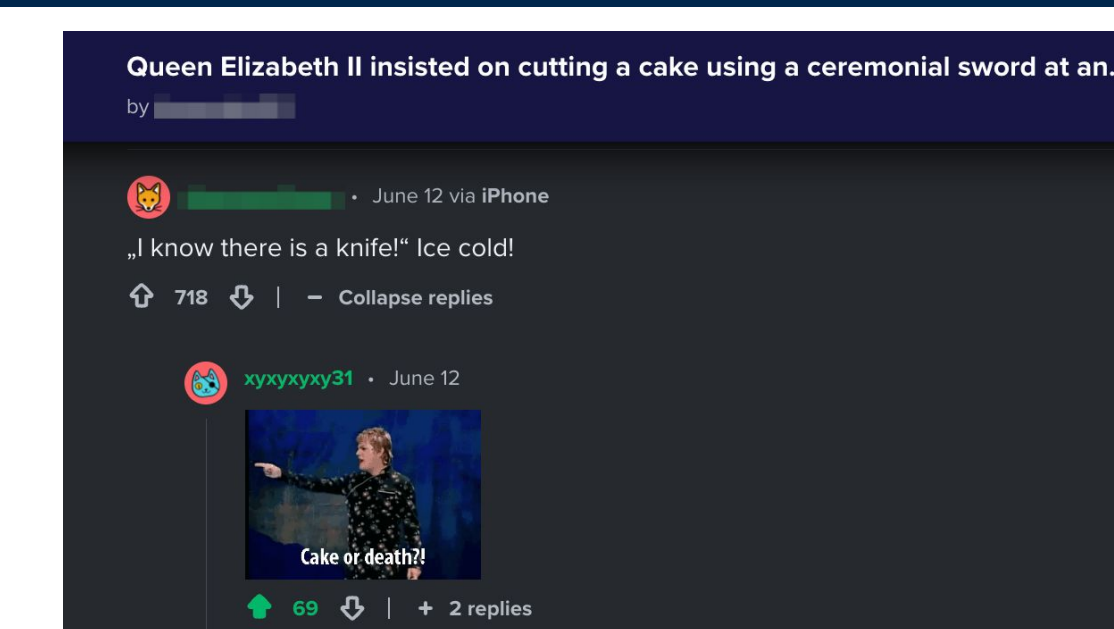
## 6. How good are our models?

How accurate were the models in picking the **exact** gif someone used in the test set?



| Model | Precision @ 1 |
|---|---|
| Select gif from predicted tags | 0.0 |
| Randomly pick a gif | 0.00002 |
| CLIP-based model | 0.000488 |
| Randomly pick a gif, relative to its frequency in the data | 0.001 |
| Pepe the King Prawn | 0.0054 |

How good are the gifs used in the replies?
(better evaluation by annotating top-10 gifs of each model)



| Model | nDCG @ 10 |
|---|---|
| Randomly pick a gif | 0.327 |
| Select gif from predicted tags | 0.453 |
| Randomly pick a gif, relative to its frequency in the data | 0.497 |
| CLIP-based model | 0.593 |
| Pepe the King Prawn | 0.815 |

(between 0 and 1, higher the better)

## 7. We deployed a Randomized Controlled Trial (RCT) in the real-world!



The RCT ran for ~5 months, and made 8,369 replies to users.

We run a regression on *score of the gif reply*, considering variables including model choice, parent comment topic, etc.



Models are compared with the effect of just randomly picking a gif

CLIP-based model ***
Randomly pick a gif, relative to its frequency in the data
Pepe the King Prawn ***
Select gif from predicted tags

Negative Binomial regression coefficients (higher is better)
(bars show standard error and *** denotes significance at 0.01)

**Code & Data:** https://github.com/xingyaoww/gif-reply
**Gif-Bot Slack App:** https://github.com/xingyaoww/gif-reply-slack-bot