



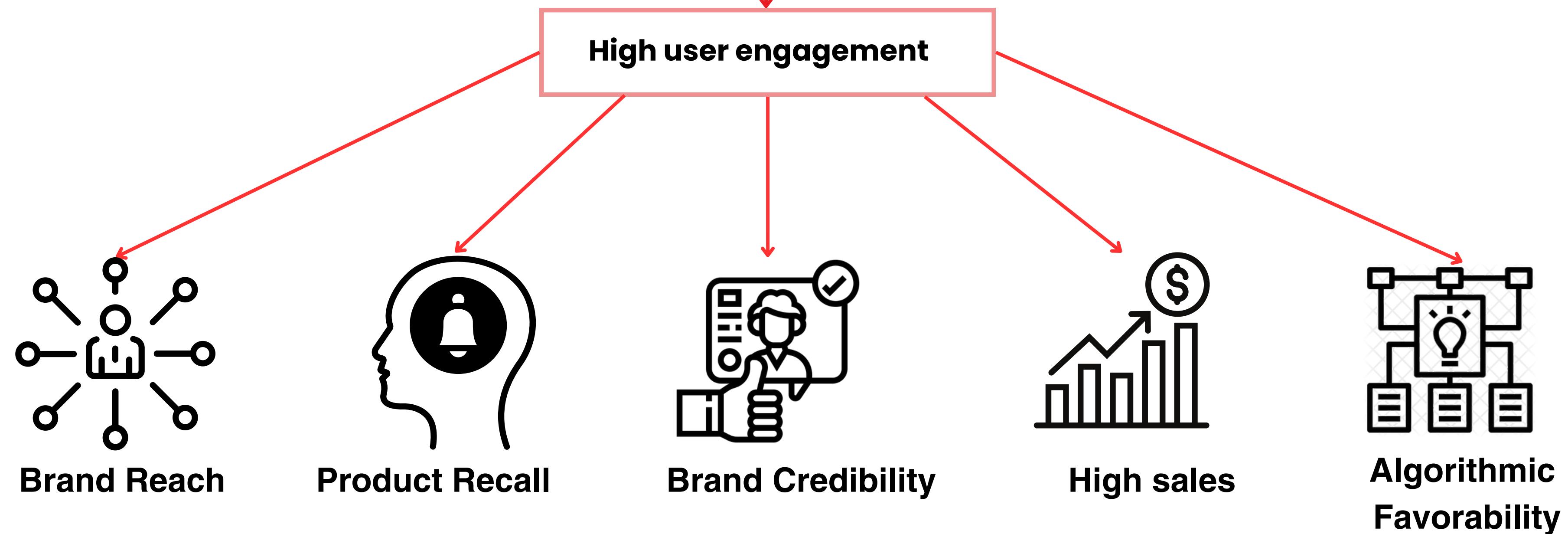
BEHAVIOUR SIMULATION CHALLENGE

TEAM 18



Background

Why is target and timely content valuable?

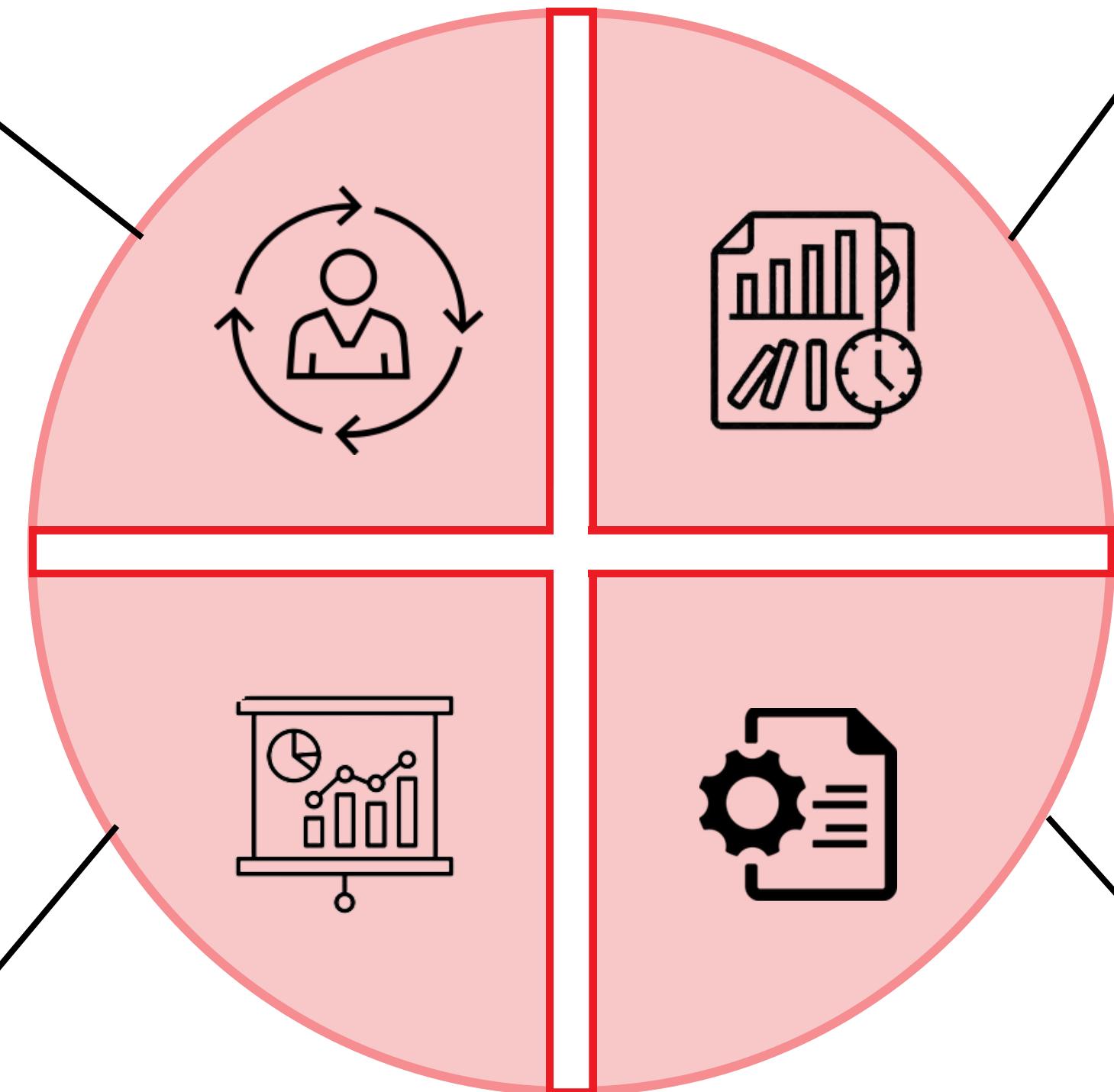




Background

PERSONALIZATION

Personalization to build strong connections between brand and consumers



PREDICTING FUTURE TRENDS

Predicting future trends to align with evolving customer needs

DATA-DRIVEN DECISIONS

Making data-driven decisions for building strategies which will yield desired outcomes

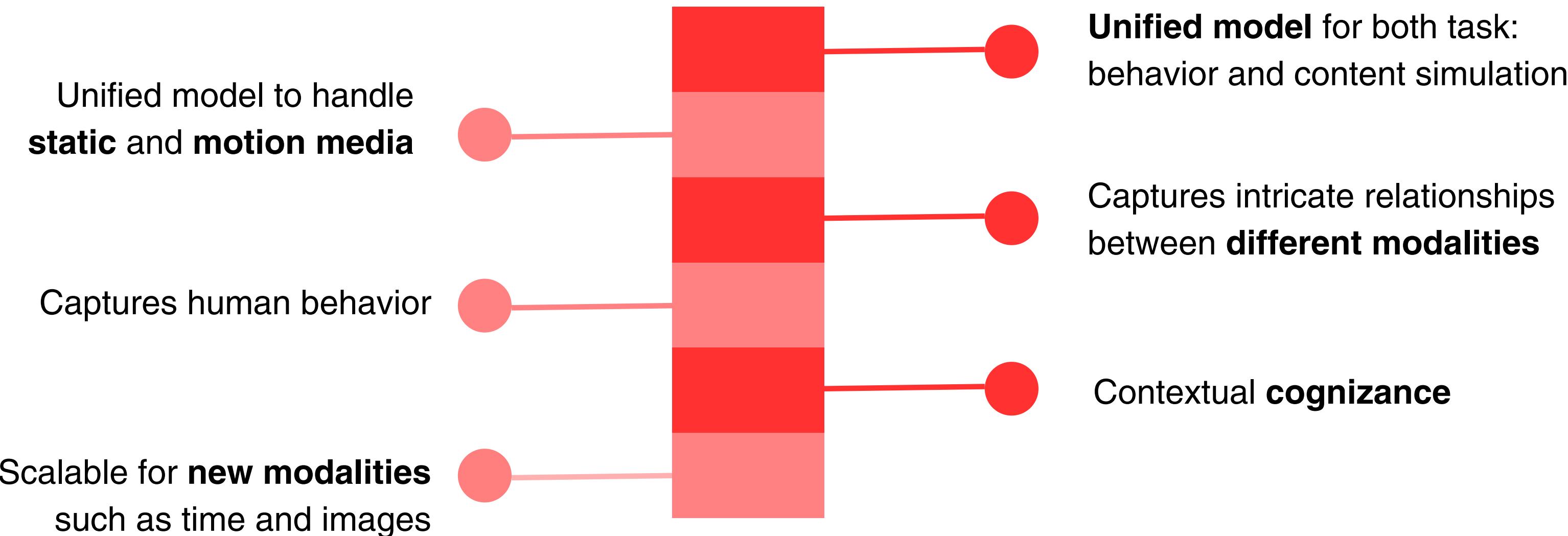
CONTENT VISIBILITY

Optimizing content visibility to increase the likelihood of content promoted

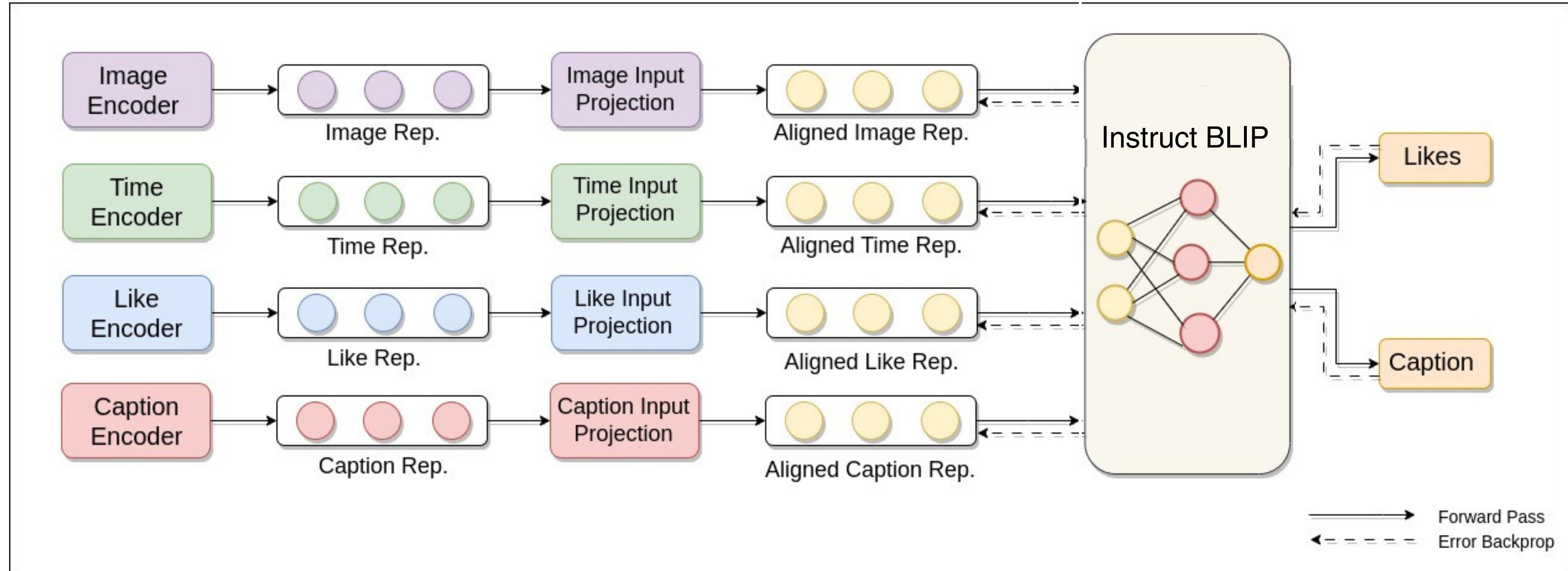


Multi-Modal Masked Learning M3L with InstructBLIP

Our Novel Approach: Multi-Modal Masked Learning (M3L)



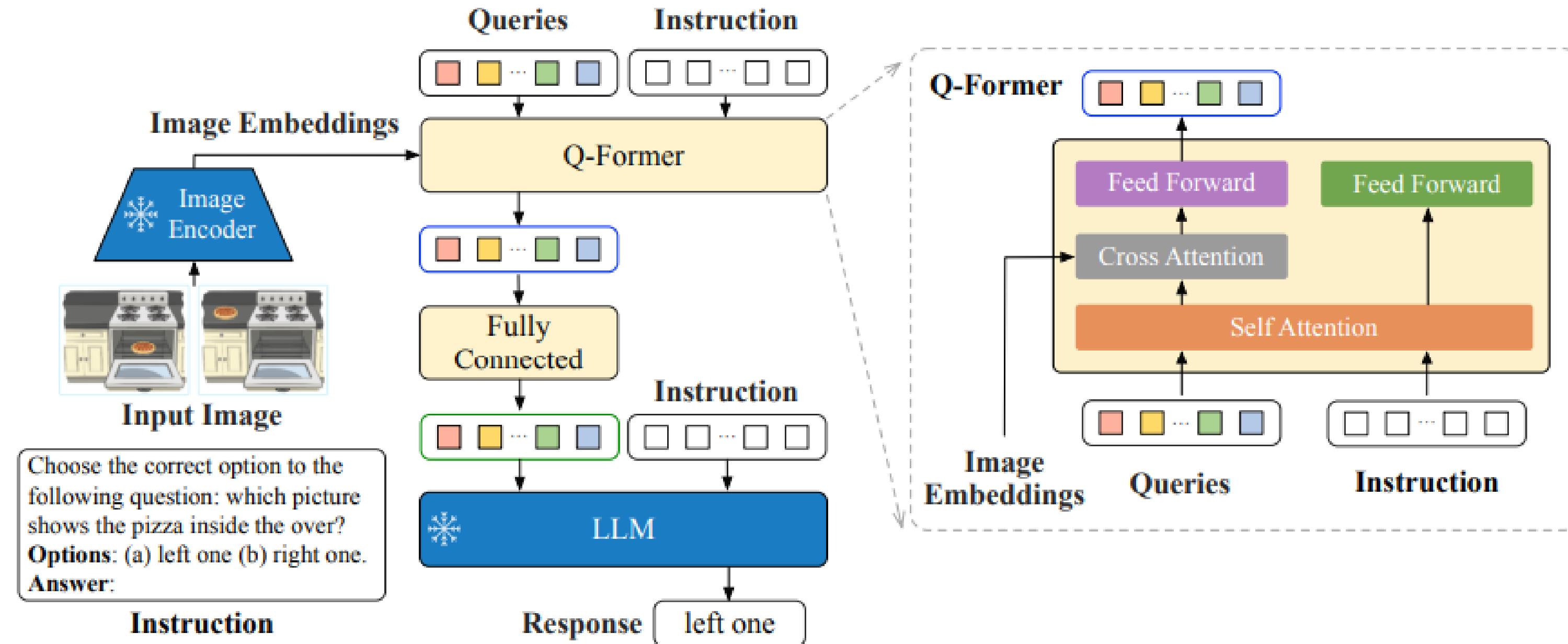
Model Pipeline



End to end model pipeline proposed for finetuning the InstructBLIP model using M3L.



Why InstructBLIP?



InstructBLIP's standout feature lies in its pioneering **Instruction-Aware Visual Feature Extraction mechanism**, setting it on a trajectory for a more insightful understanding of user interactions and responses in social media.

[arXiv:2305.06500 \[cs.CV\]](https://arxiv.org/abs/2305.06500)



Video Frame Retriever



Duration :60 seconds

Frame rate : 30 frames per second

Total media volume = **1800 frames**

An 85-year-old primary school in Shanghai has been lifted off the ground — in its entirety — and relocated using new technology dubbed the "walking machine."



The 85-year-old primary school was lifted from the ground and relocated about 200 feet



A school in Shanghai was rotated and moved using "walking machine" technology



This is the first time Shanghai has used the "walking machine" method to relocate a building



Yes, that building is walking



The 85-year-old primary school was lifted from the ground and relocated about 200 feet

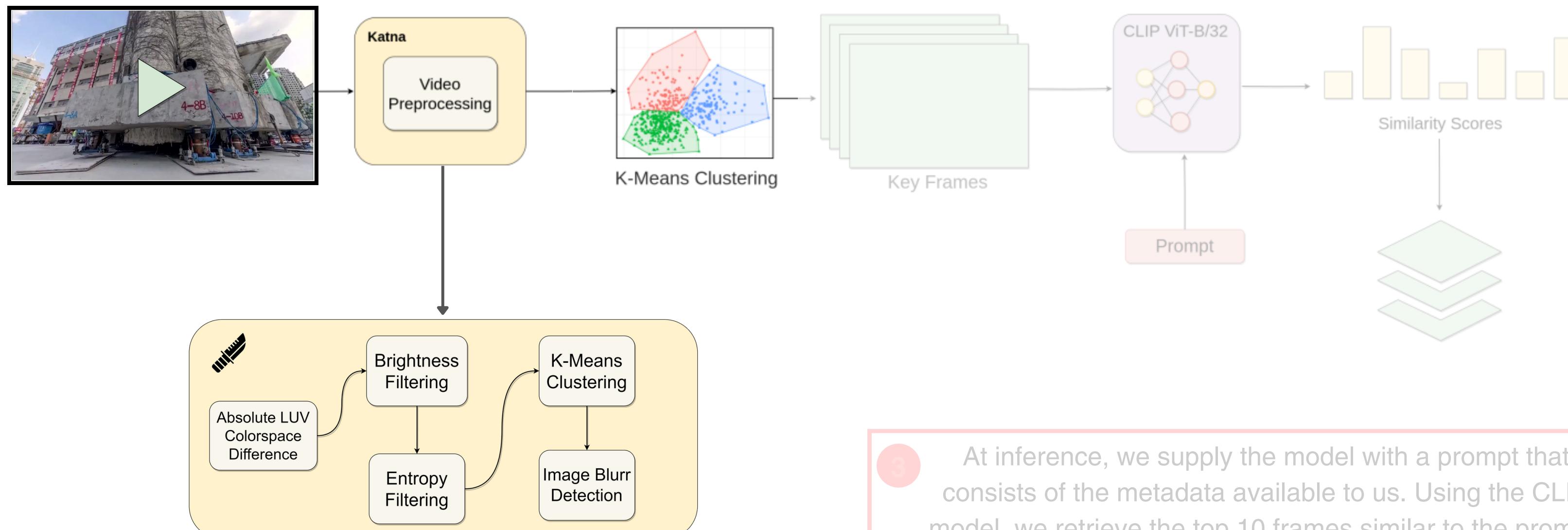
Sufficient to provide contextual information !



Video Frame Retriever

1 The pipeline consists of a **two-step distillation process**, that includes sampling frames at regular intervals

2 Then, we leverage the Katna library to extract the top **50 keyframes** from the sampled frames.

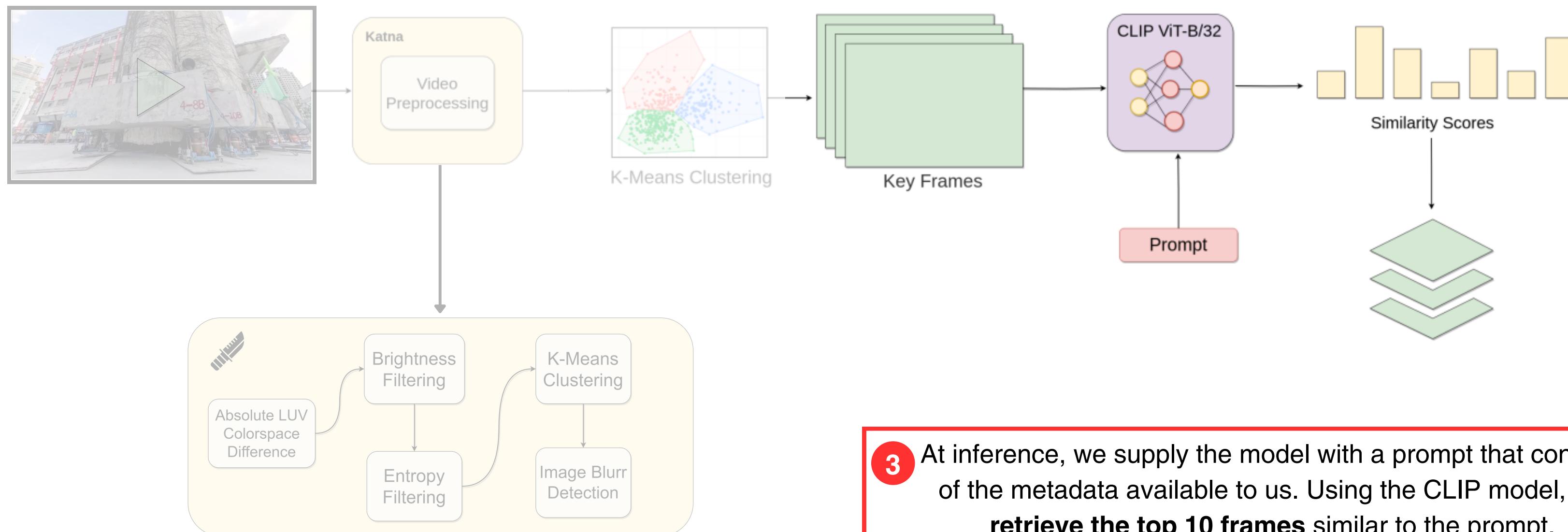




Video Frame Retriever

1 The pipeline consists of a two-step distillation process, that includes sampling frames at regular intervals

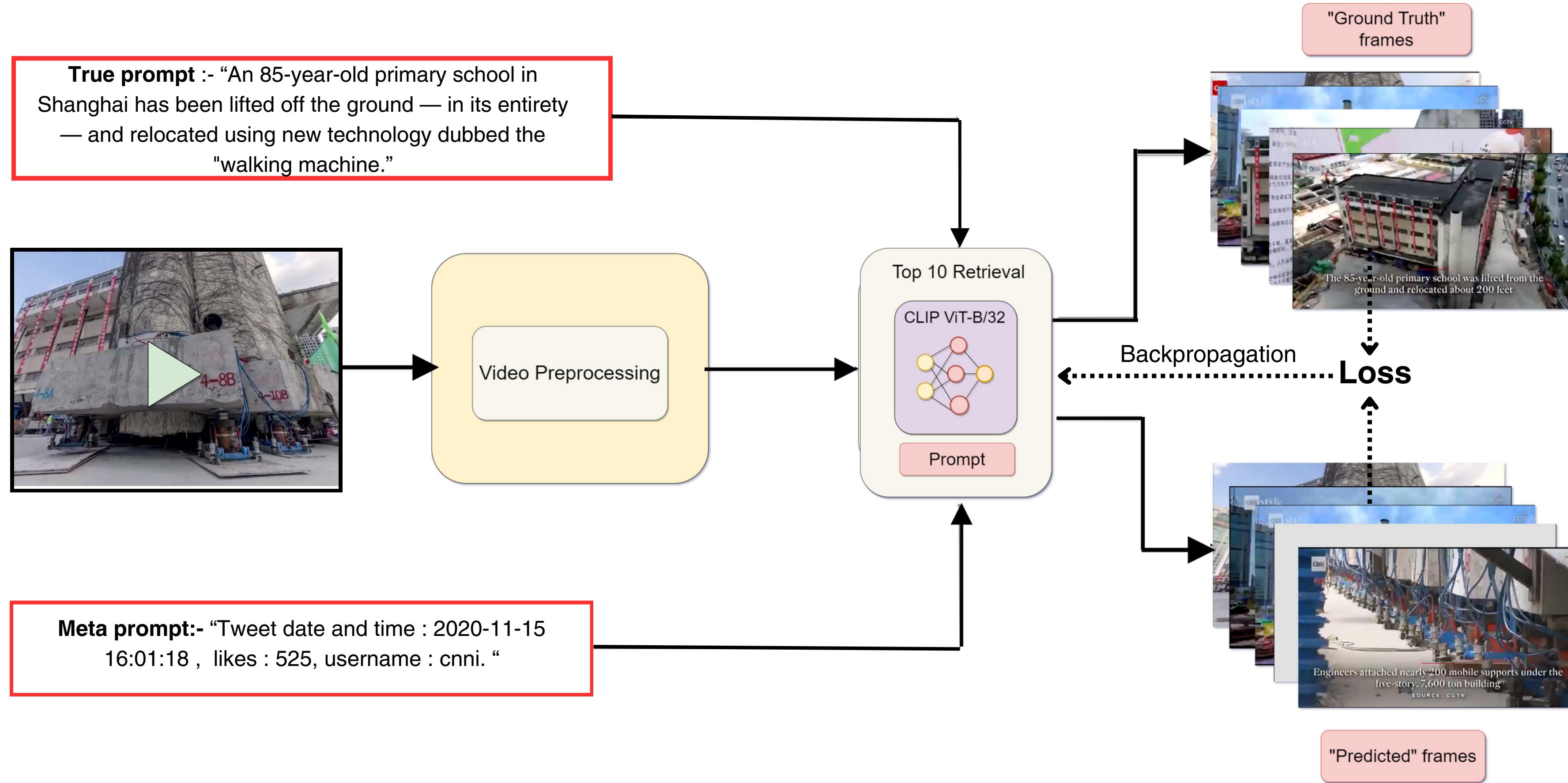
2 Then, we leverage the Katna library to extract the top 50 keyframes from the sampled frames.



3 At inference, we supply the model with a prompt that consists of the metadata available to us. Using the CLIP model, we **retrieve the top 10 frames** similar to the prompt.



Training Procedure



Training Procedure

Custom Loss Function L for optimising model parameters

$$DFRloss(q_i, fp_i^+, fn_{i,1}^-, \dots, fn_{i,n}^-) = -\log \frac{e^{\text{sim}(q_i, fp_i^+)}}{e^{\text{sim}(q_i, fp_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, fn_{i,j}^-)}}$$

$$L = 0.1 \times BCEloss + 0.9 \times DFRloss$$

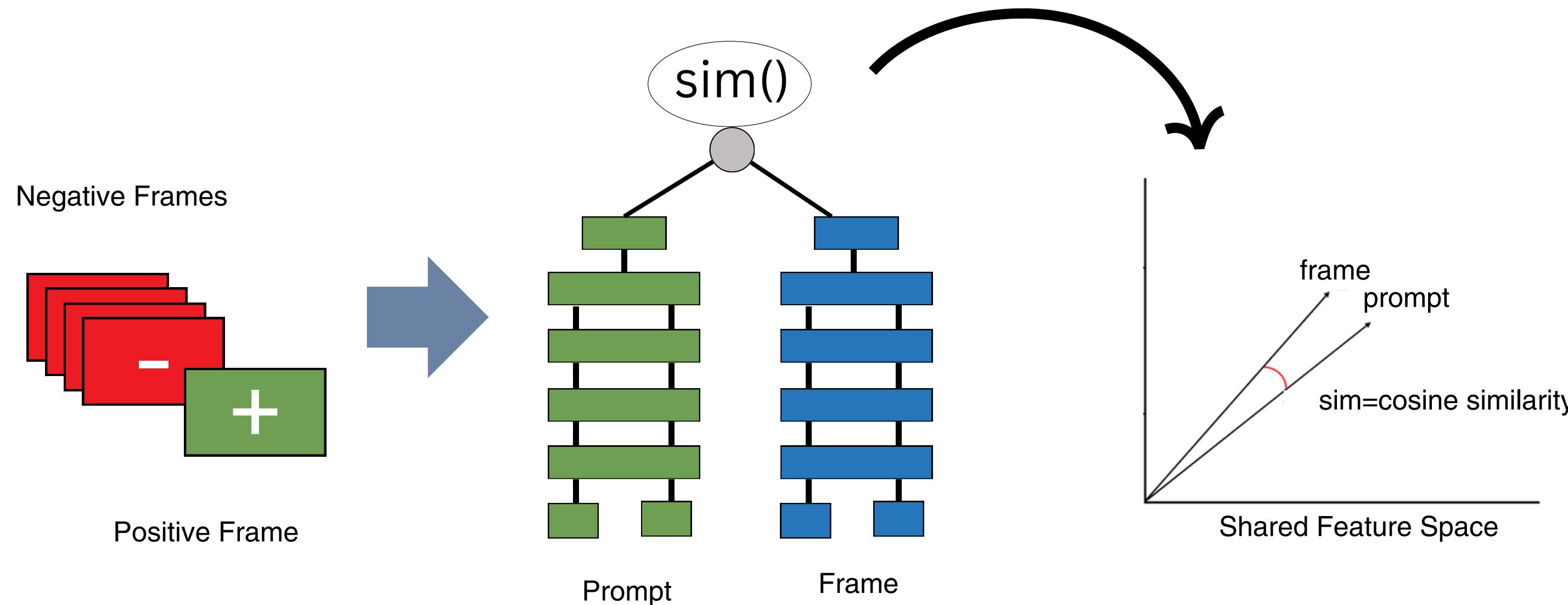
$$BCEloss = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$



Training Procedure

$$DFRloss(q_i, fp_i^+, fn_{i,1}^-, \dots, fn_{i,n}^-) = -\log \frac{e^{\text{sim}(q_i, fp_i^+)}}{e^{\text{sim}(q_i, fp_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, fn_{i,j}^-)}}$$

Inspired by the loss function presented in "**Dense Passage Retrieval for Open-Domain Question Answering**"



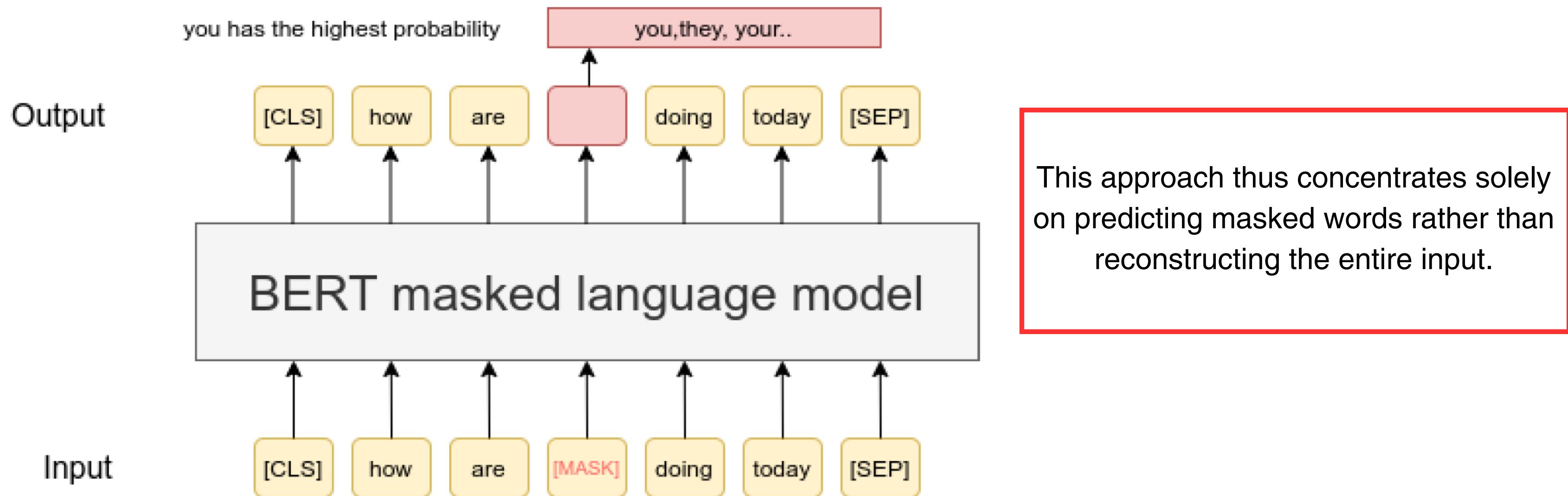
[arXiv:2004.04906 \[cs.CL\]](https://arxiv.org/abs/2004.04906)



Masked Language Modelling

The "Masked Language Modelling" (MLM) technique involves randomly **masking** a small percentage of input tokens in sequences to predict these hidden tokens by feeding their final representations into an output softmax.

3



* arXiv:1810.04805



Training InstructBLIP using Multi-Modal Masked Learning

We carry out MLM on Multimodal data in our M3L approach for fine-tuning the InstructBLIP model for content and behaviour simulation to enhance the model's understanding of human behavior by providing it context-aware representations by predicting masked elements within the text.

Prompt for Content Simulation using M3L

You are content creating wizard. Today I challenge you to predict a caption for an post that was posted in the past. I'll provide you with the following information regarding the post:

Username, Date of post release, Time of post release and No of likes

Review the username and date as there could be some insights relevant for the caption. Identify the material/person/objects in the image and align the caption accordingly. This will help to predict a better caption.

Now given the image,

Username is MayoClinic.,

Release Date is 2020-07-19,

Release Time is 00:07:00,

Likes received is 78.

Predicted Caption: [MASK]



Training InstructBLIP using Multi-Modal Masked Learning

In the prompt for behavior simulation we mask the captions modality and ask the model to predict the same given the other modalities as depicted in the prompt below.

Prompt for Behavior Simulation task using M3L

You are content creating wizard. Today I challenge you to predict a the number of likes for a post that was posted in the past. I'll provide you with the following information regarding the post:

Username Date of post release Time of post release Caption

Review the username and date as there could be some insights relevant for the number of likes. Identify the material/person/objects in the image and obtain the number of likes. . This will help to predict the number of likes.

Now given the image,

Username is MayoClinic.,

Release Date is 2020-07-19,

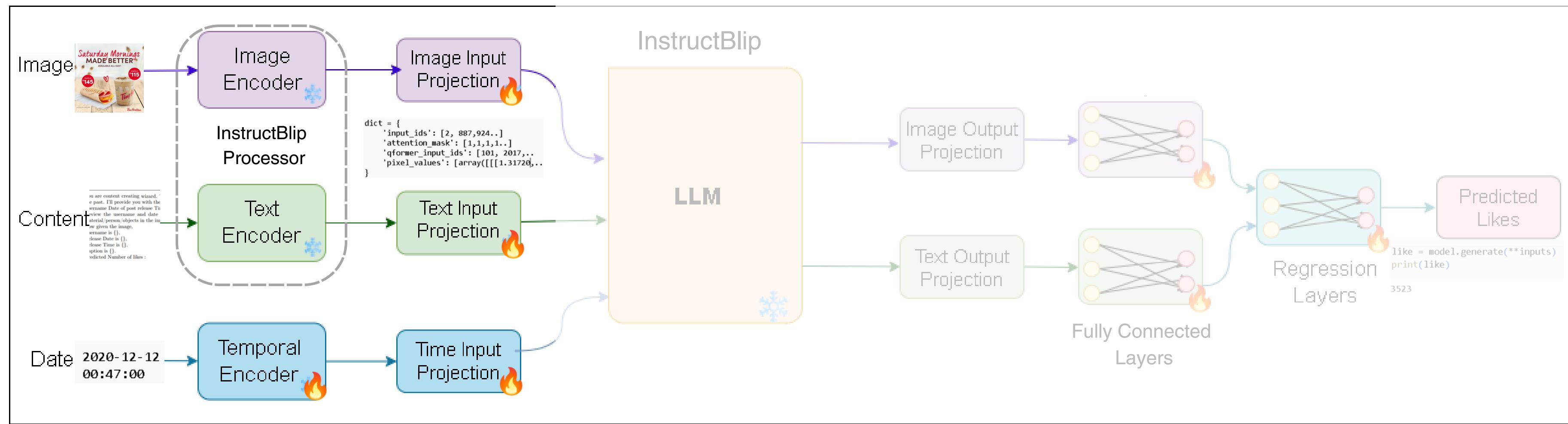
Release Time is 00:07:00,

Caption is Many people are understandably anxious to get outdoors in light of quarantines. Do it safely with these tips for outdoor activities: <hyperlink> COVID19 <hyperlink>.

Predicted Number of likes : [MASK]

Training for the Behaviour Simulation Task

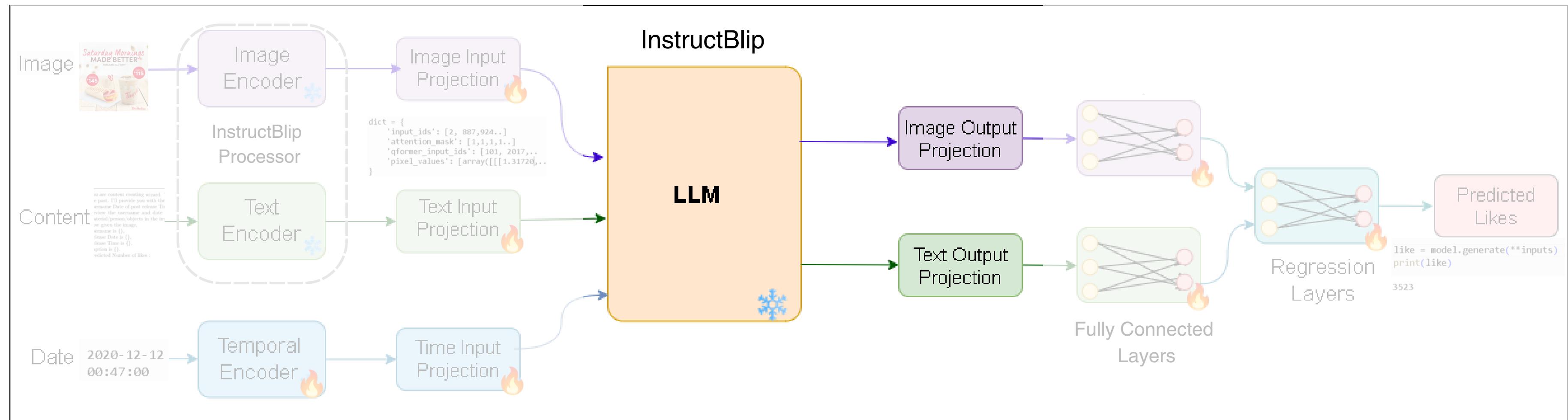
- 1 The text modality and image modality is processed by InstructBlip processor which consists of
 - a. BLIP as image encoder and
 - b. Llama/t-5 tokenizer as text encoder.





Training for the Behaviour Simulation Task

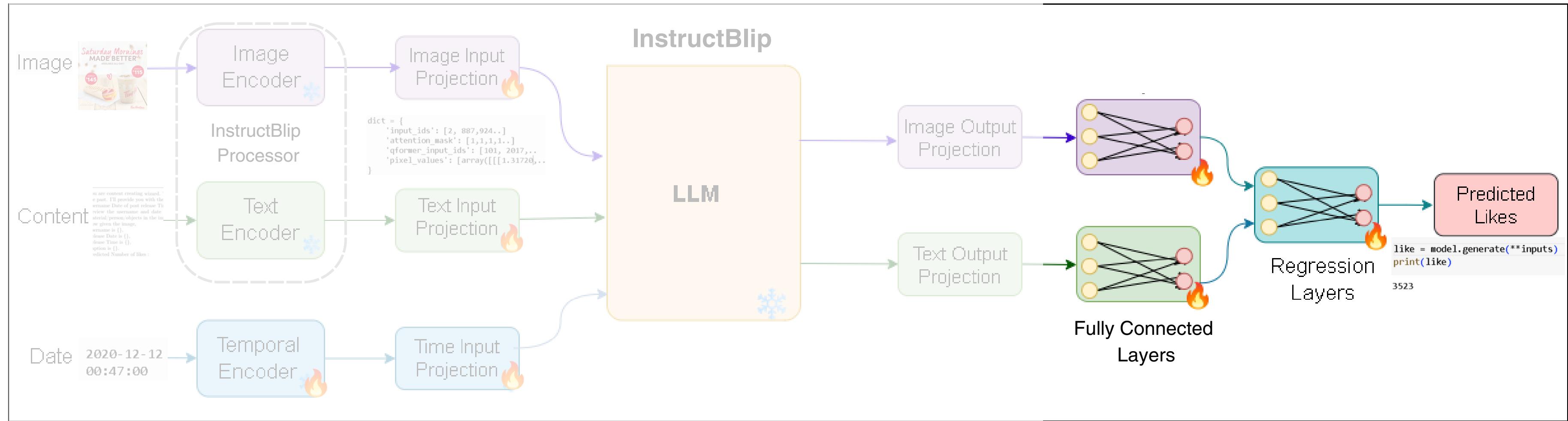
- ② From the output of the model, last hidden states is extracted for both text and image encodings





Training for the Behaviour Simulation Task

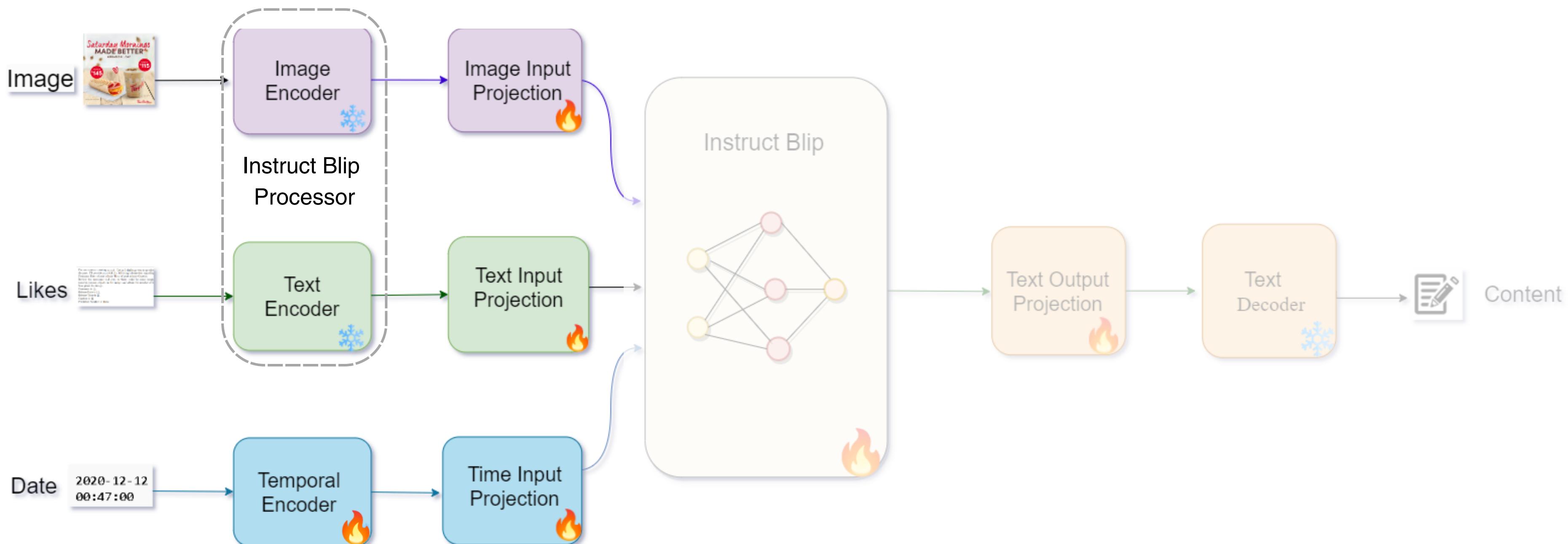
These encodings are passed
③ through various connected and
regressor layers before
prediction of final output.





Training for the Content Simulation Task

- 1 The text modality and image modality is processed by InstructBlip processor which consists of
- BLIP as image encoder and
 - Llama/t-5 tokenizer as text encoder.

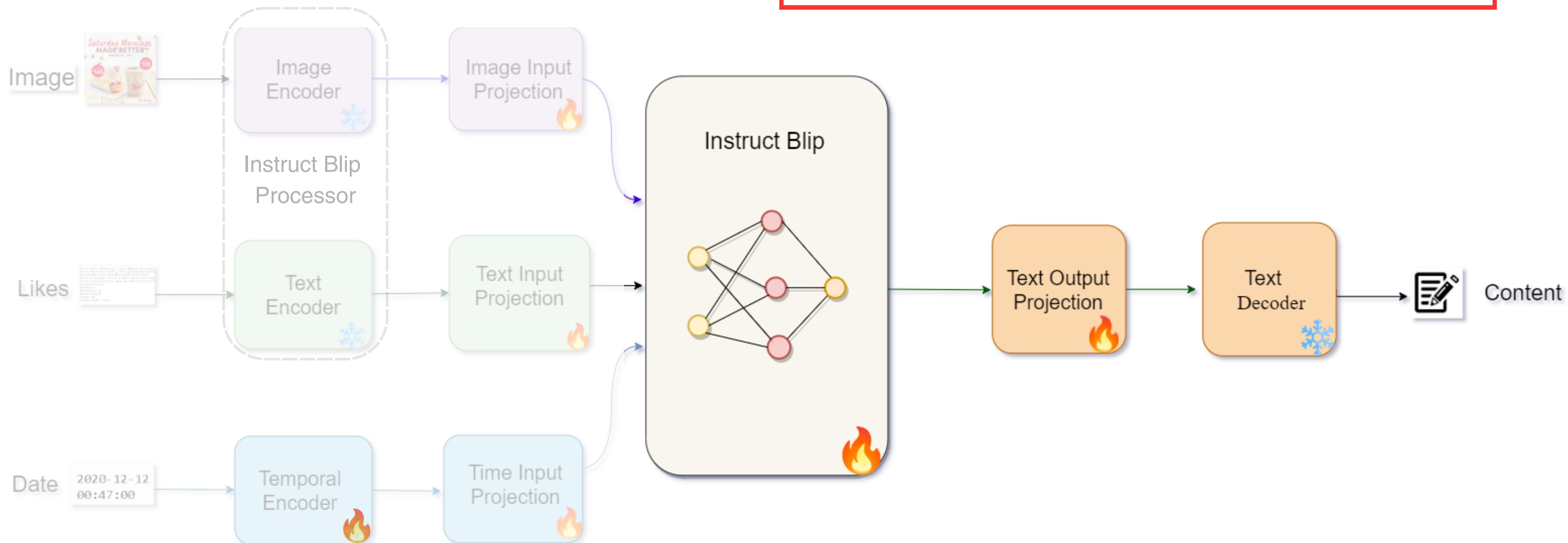




Training for the Content Simulation Task

2

We derive outputs from the language output head of the InstructBLIP Vicuna model and pass it through a batch decoder to get the predicted captions





Training Setup

Algorithm 1 Masked Vision Language Modelling for Behaviour Simulation

```
1: InstructBLIP Model (M):
2:   M  $\leftarrow$  InstructBLIP Model
3:   M_Q  $\leftarrow$  Q-Former part of M
4:   M_T  $\leftarrow$  T-Former part of M processes temporal information.
5:
6: Initialization:
7: TWEET_METADATA  $\leftarrow$  Extract tweet metadata from  $D_{train}$ 
8: LIKES_PROMPT  $\leftarrow$  Initialize likes_prompt
9: CAPTIONS_PROMPT  $\leftarrow$  Initialize captions_prompt
10: SHARED_FEATURES  $\leftarrow$  Propagate batches through shared base layers
11: T_FORMER  $\leftarrow$  Create Temporal-Former based on Q_FORMER
12:
13: Threshold Constraint:  $0 \leq THRESHOLD \leq 100$ 
14:
15: procedure MASKEDVISIONLM( $D_{train}, EPOCHS$ )
16:   for epoch  $\leftarrow 1$  to  $EPOCHS$  do
17:     for each batch in  $D_{train}$  do
18:       BATCH_CONTENT  $\leftarrow$  LIKES_PROMPT
19:       BATCH_BEHAVIOR  $\leftarrow$  CAPTIONS_PROMPT
20:
21:       if THRESHOLD  $> 50$  then
22:         LOSS  $\leftarrow$  Negative log likelihood loss for caption generation
23:       else
24:         LOSS  $\leftarrow$  MSE loss for likes generation
25:       end if
26:
27:       Backpropagate on the loss and update model parameters  $\theta$ 
28:     end for
29:   end for
30: end procedure
```

$$MSEloss = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$NLLloss = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i)]$$



How does InstructBLIP learn to model behaviour?

Training for
Behaviour Simulation

Training for
Content Simulation



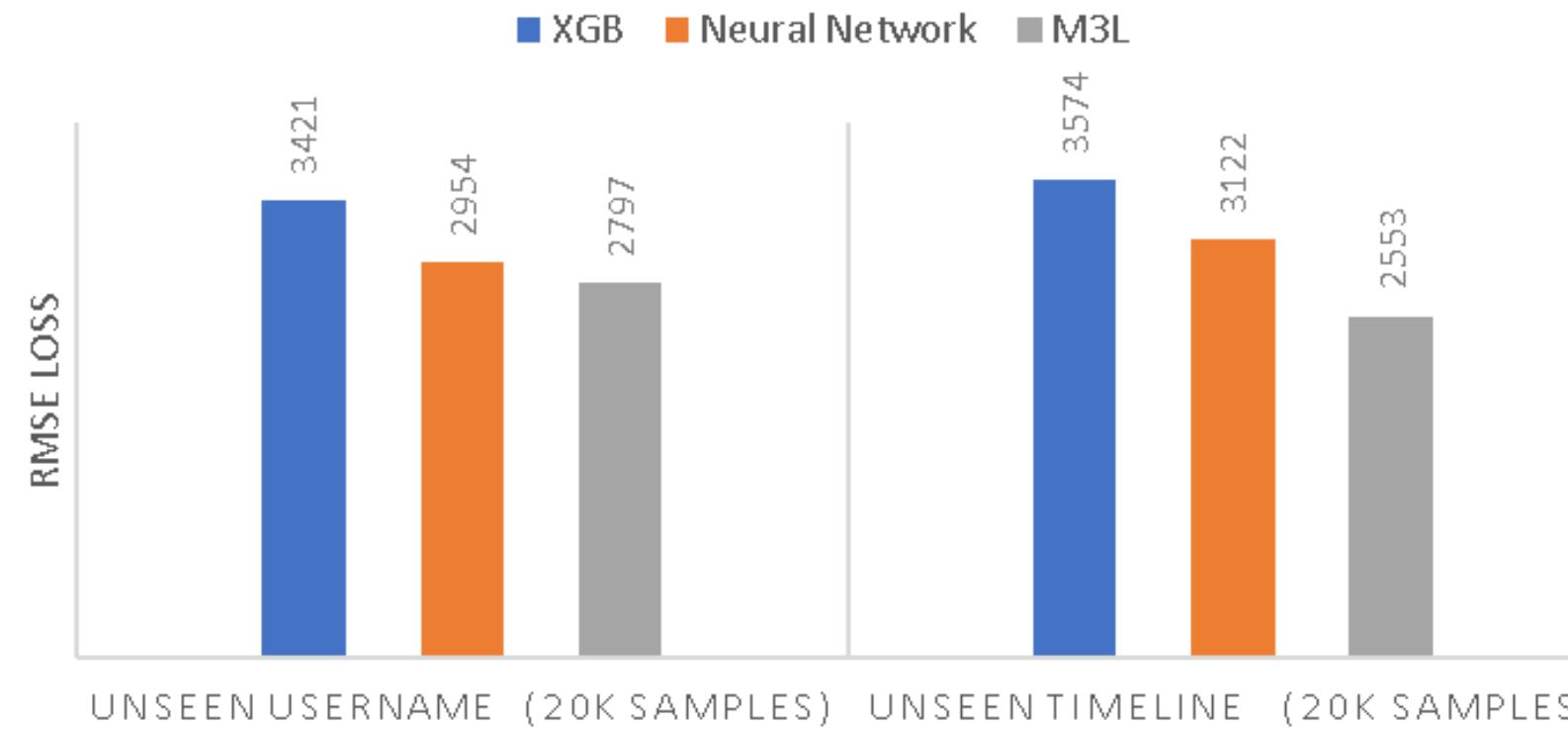
Behaviour Optimized
InstructBLIP

The multitasking fine-tuning with a custom loss in InstructBLIP refines caption generation aligned with human preferences and predicts engagement through nuanced cues.

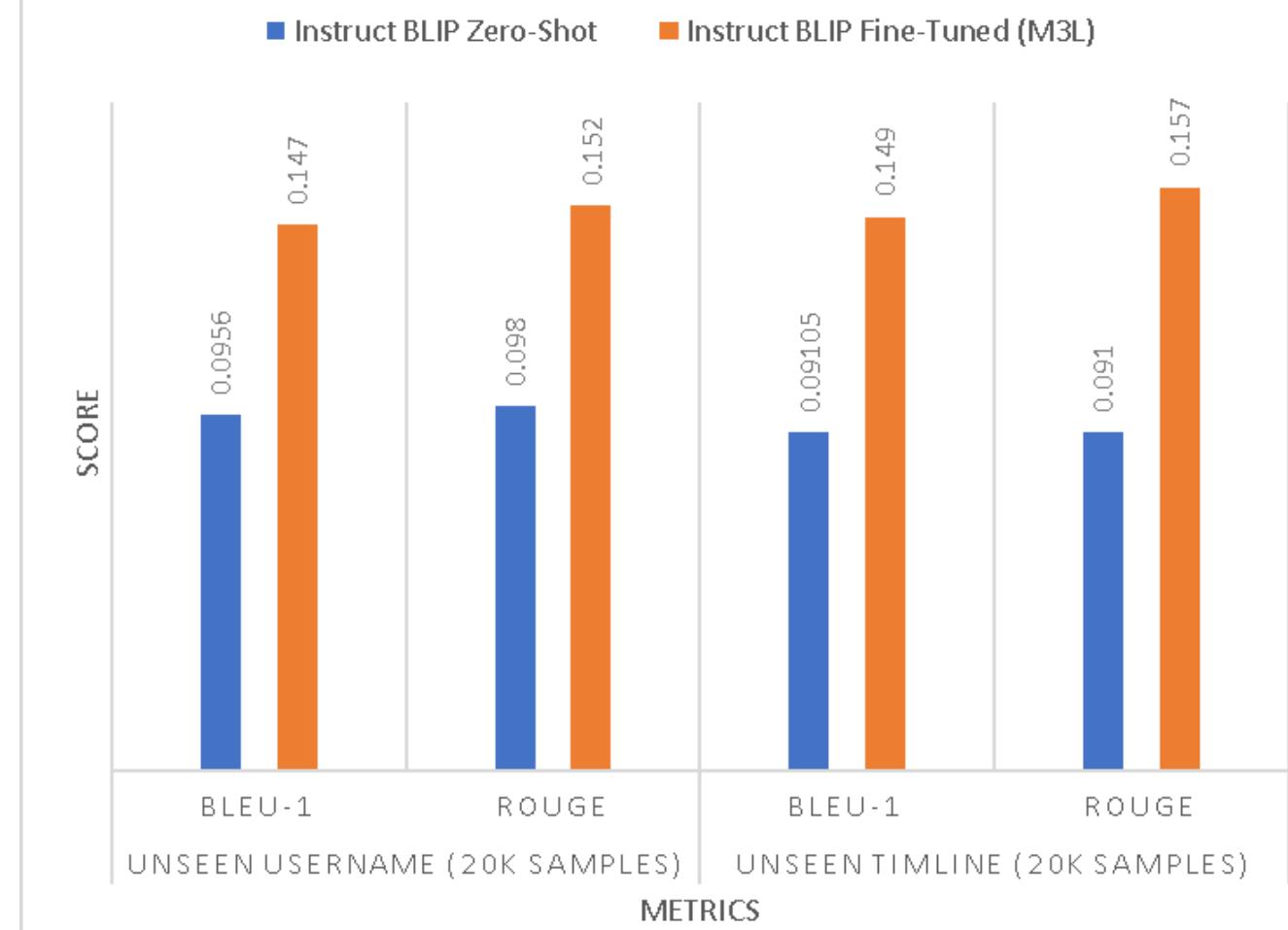


Model Performance

Behaviour Simulation



Content Simulation



*Results evaluated on training data provided by Adobe for the Behaviour Simulation Challenge



Demo Video of Inference

We can further extend our model and train it for multiple tasks to help it inculcate a deeper knowledge of human behaviour and equip it to predict more output modalities.