

3rd
13th
MVIP2024



MVIP2024



13th Iranian and 3rd International Conference
on Machine Vision and Image Processing

6 & 7 March , 2024

MVIP2024.KHU.AC.IR

Cultural-aware AI model for emotion recognition

Mehrdad Baradaran

Shahid Beheshti University

Payam Zohari

Khajeh Nasir Nniversity of Technology

Abtin Mahyar

Institute of Research in Fundamentals

Hossein Motamednia

Institute of Research in Fundamentals

Dara Rahmati

Shahid Beheshti University

Saeid Gorgin

Chosun University

Presentation Outline:

- Problem definition
- Dataset introduction
- Proposed methods
 - Image Model
 - Text Model
- Linear Combination Model
- Results
- Conclusion

Do cultures affect our interpretation?



شلال طبيعي جميل. مشاعر النمو والحيوية والطاقة موجودة.

Translation: Beautiful natural waterfall. Feelings of growth, vitality and energy.

Excitement
Arabic



The water that's rushing downward looks like a bride's wedding veil.

Awe
English



瀑布就像四蹄生风的白马如潮水涌来，非常的壮观

Translation: The waterfall is like a white horse and wind, it is spectacular.

Contentment
Chinese



How to utilize captions to Convey emotions?

WIKIART

Construction of an emotional image captioning model, requires the task of image processing as well as emotion recognition. The “WIKIART” dataset, supplying thousands of classified paintings helped us through this journey.



ArtELingo

To approach cultural-awareness model, we needed to train our model on a dataset rich in multilingual captions and “ArtELingo” provided us with a collection of 80k annotated artworks in 3 languages, namely English, Arabic, and Chinese.

But why “ArtELingo”?

	COCO	ArtEmis	ArtELingo
Image Source	Photos	WikiArt	WikiArt
#Images	328k	80k	80k
#Annotations	2.5M	0.45M	1.2M
#Annot/Image	7.6	5.68	15.3
Emotions	0	9	9
Languages	E	E	ACES

Proposed Method

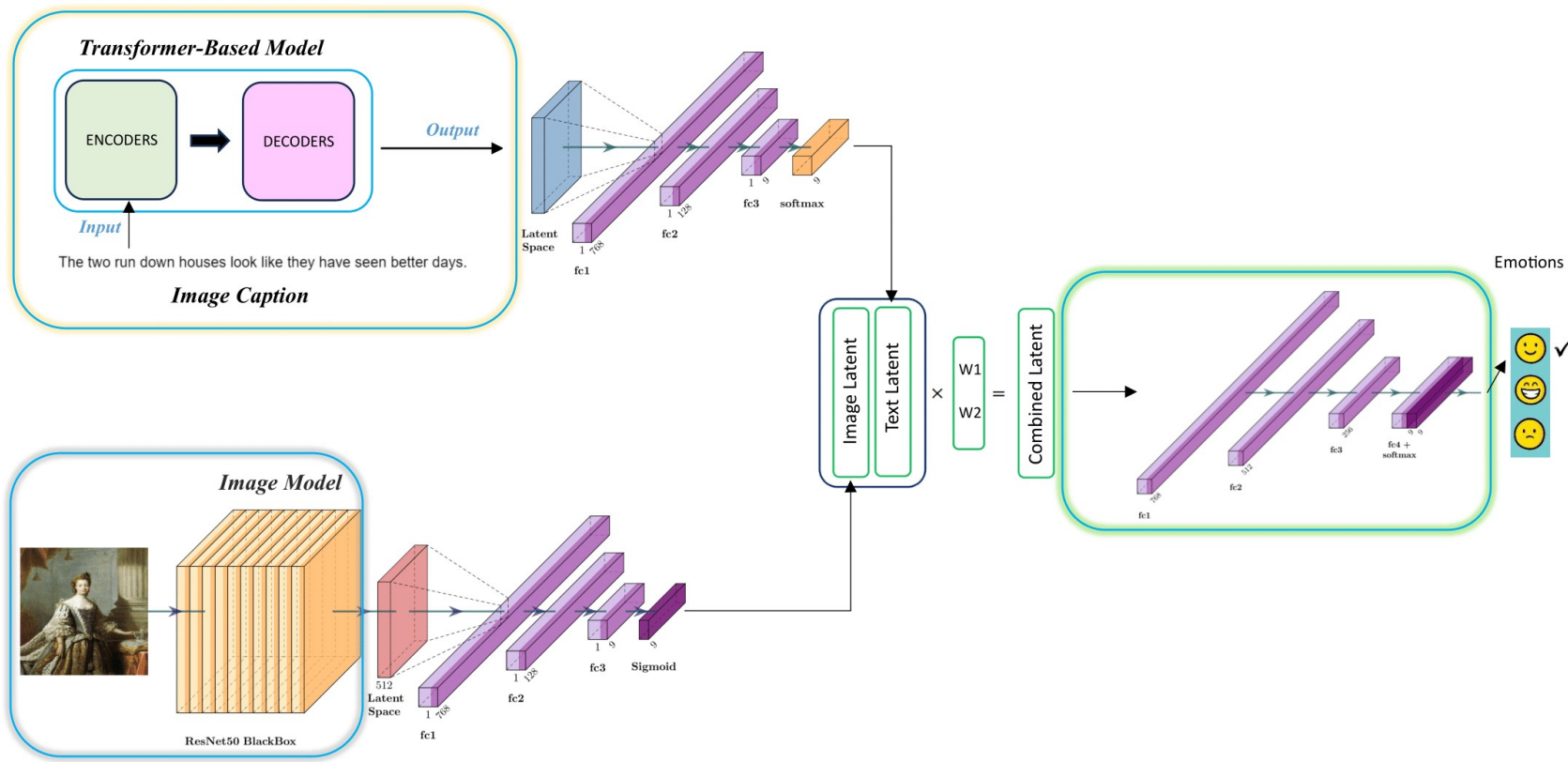
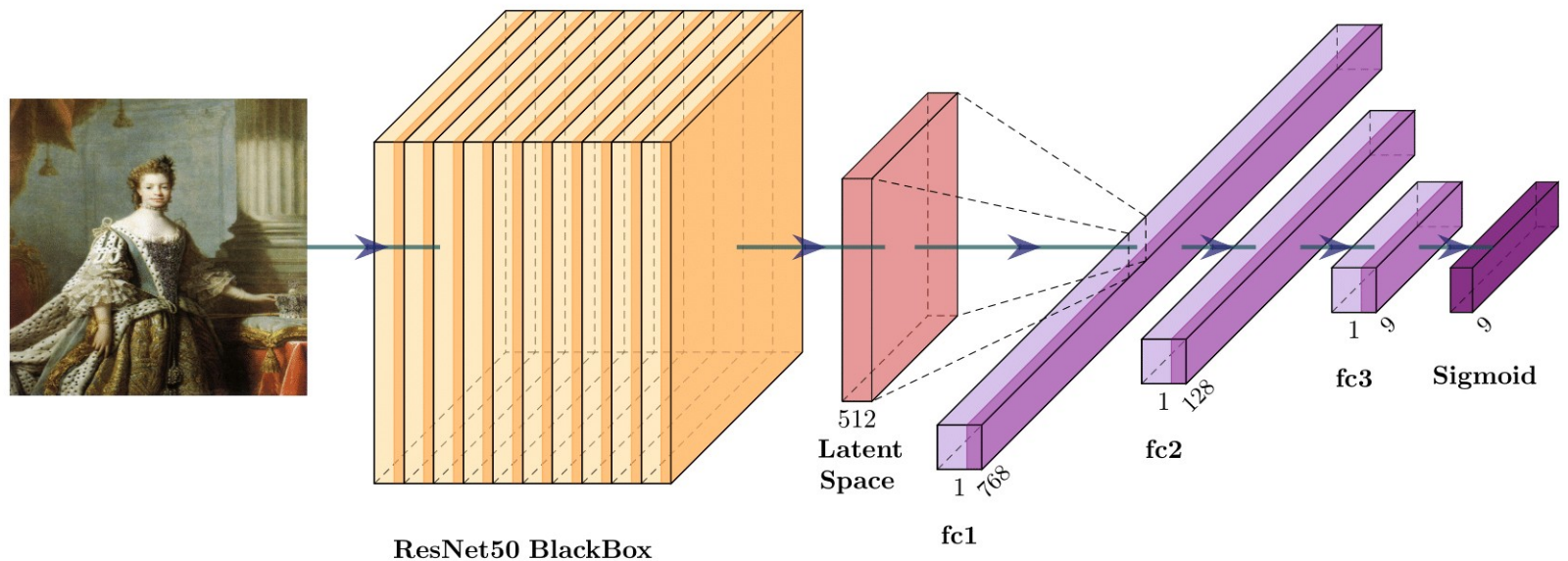
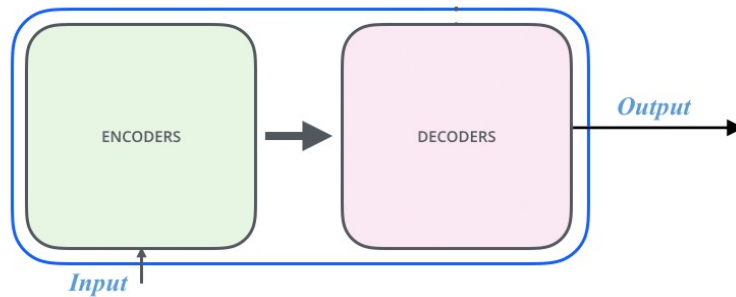


Image Model



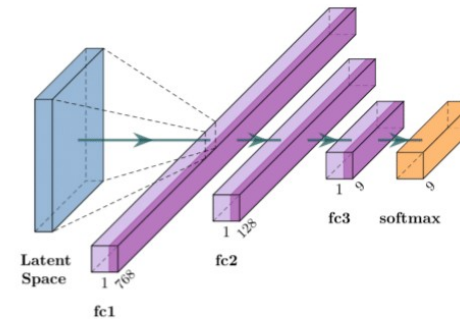
Text Model

Transformer-Based Model

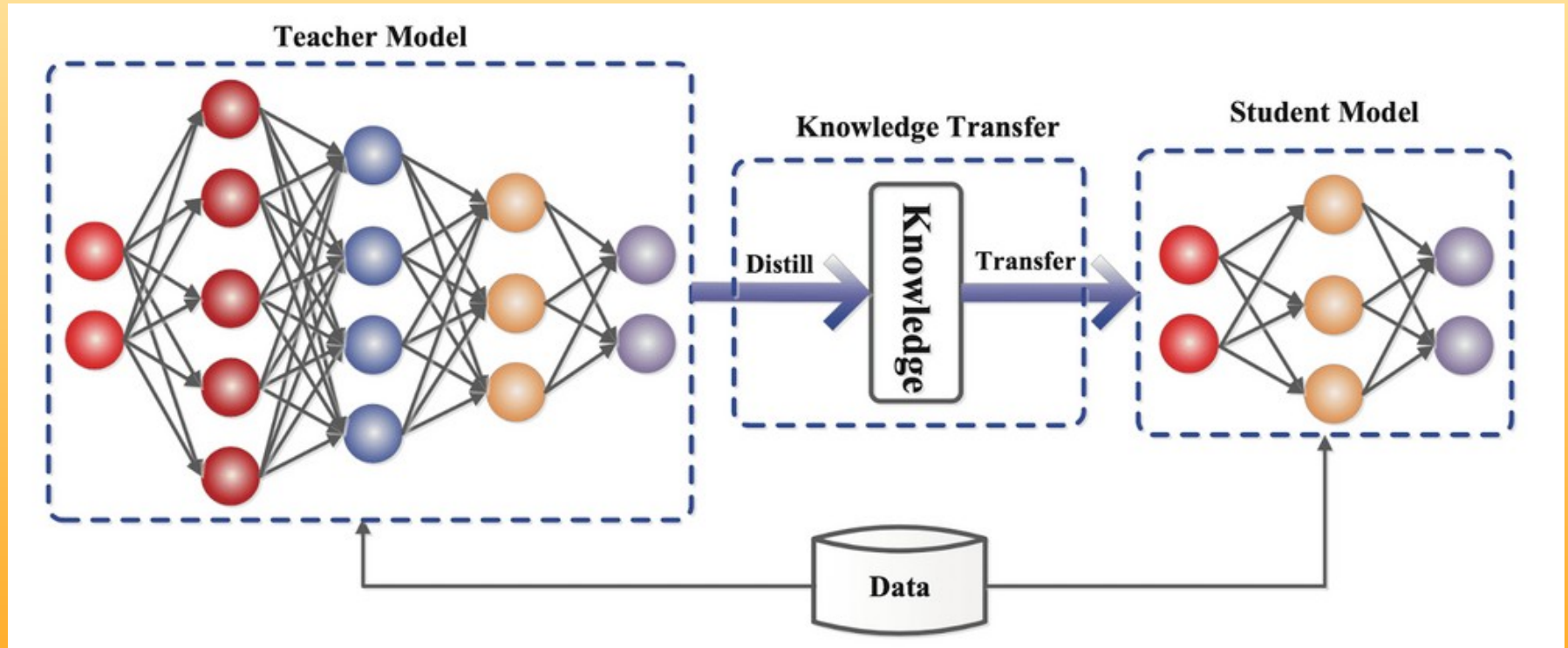


The two run down houses look like they have seen better days.

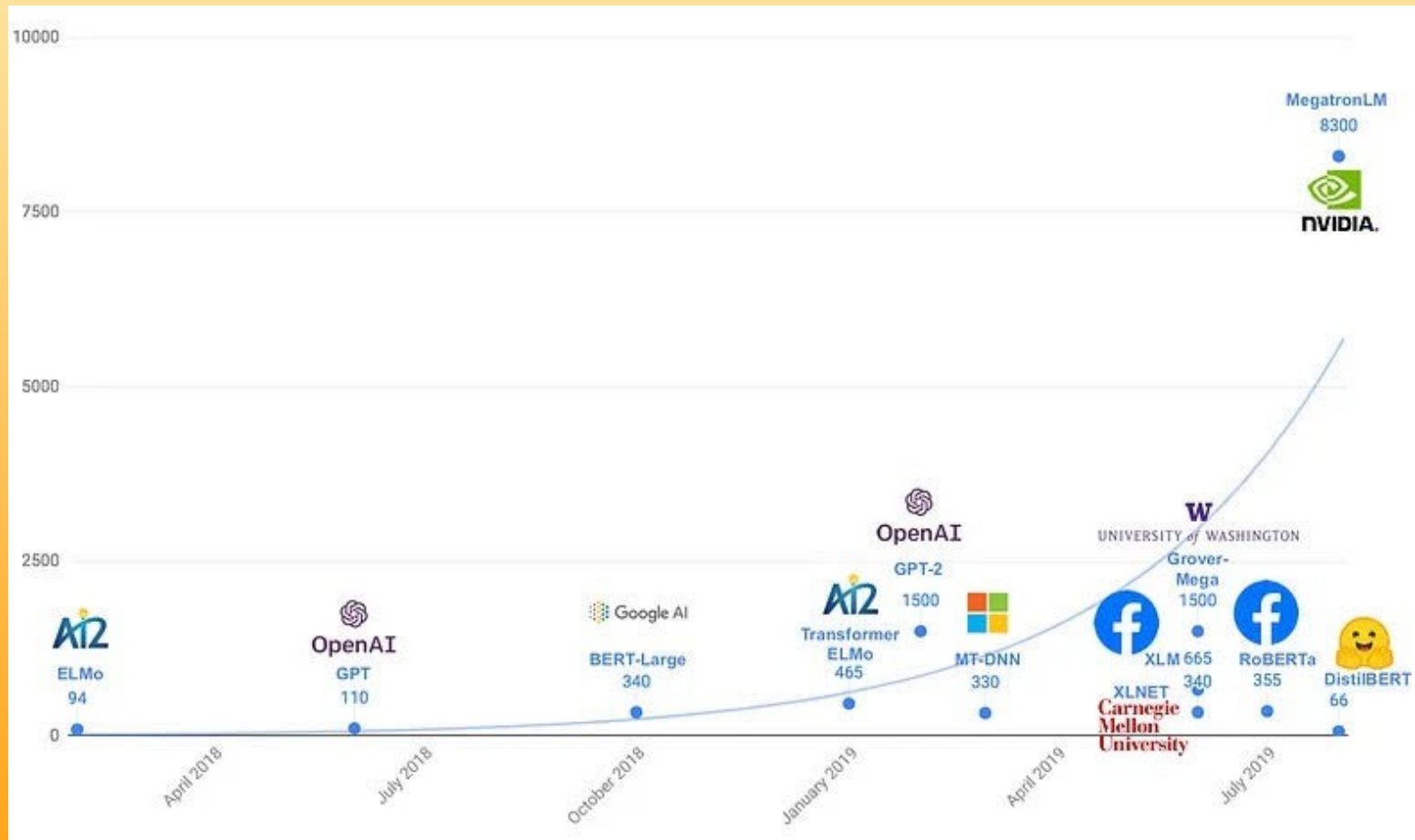
Image Caption



Knowledge Distillation



Why DistilBERT?



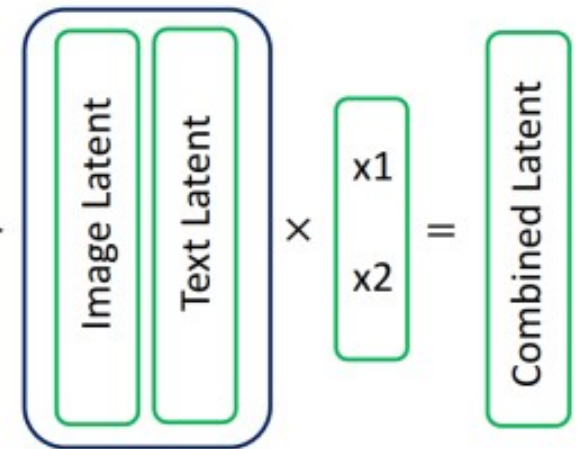
Feature Linear Combination

Linear Combination of **Image Latent** and **Text Latent**

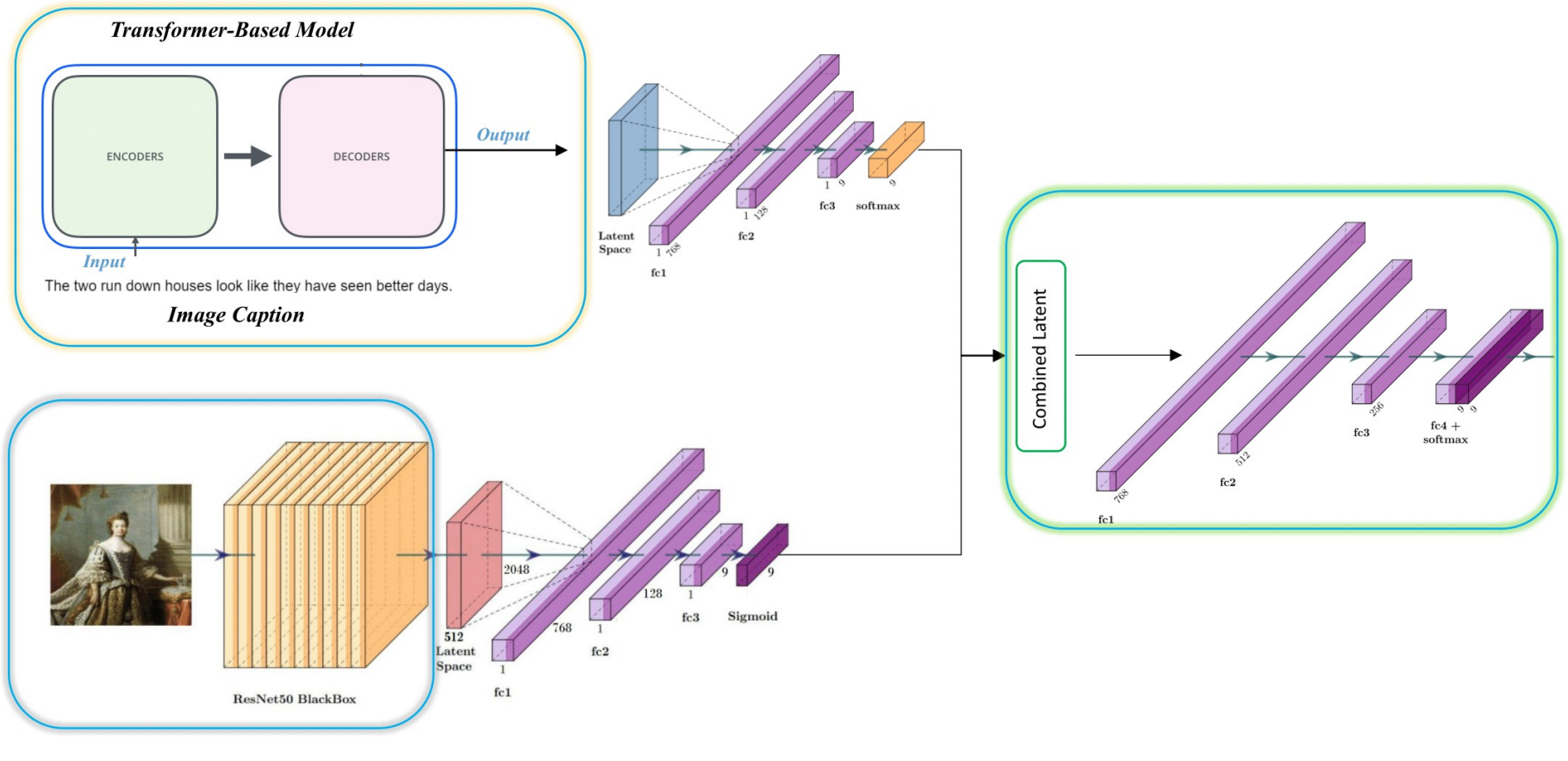
$$x1 \cdot \text{Image_Latent} + x2 \cdot \text{Text_Latent}$$

Scalars

$$Ax = b$$



Model Overview

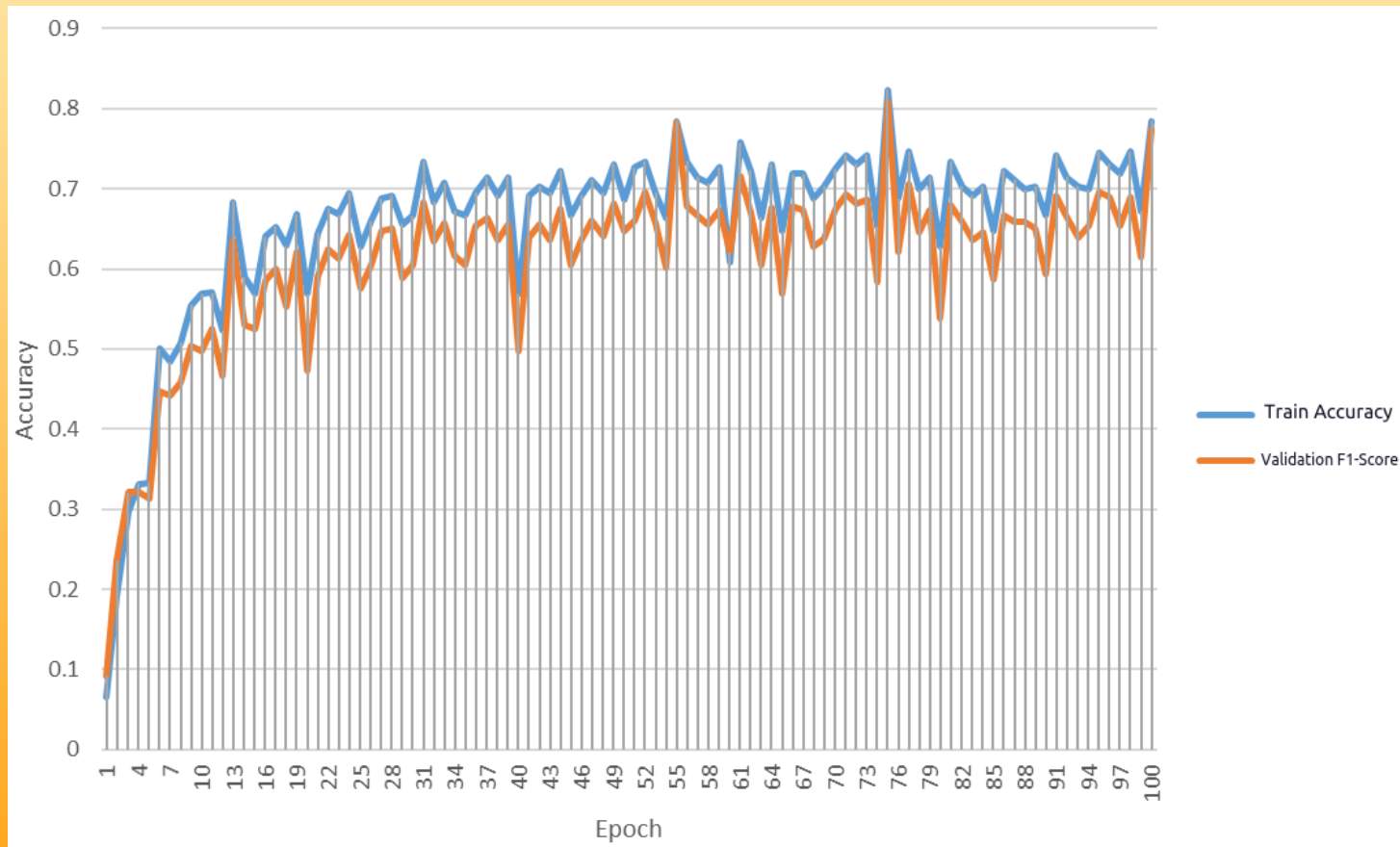


Our Results

THE EVALUATION OF THE IMAGE MODELS ON THE ARTELINGO DATASET

Model/Accuracy	Train ACC	Validation Acc
ResNet50	0.711386	0.698935688
VGG	0.717367	0.706483998
VIT	0.702534	0.687801932

Our Results



Thanks for your attention!