

HATEFUL MEMES

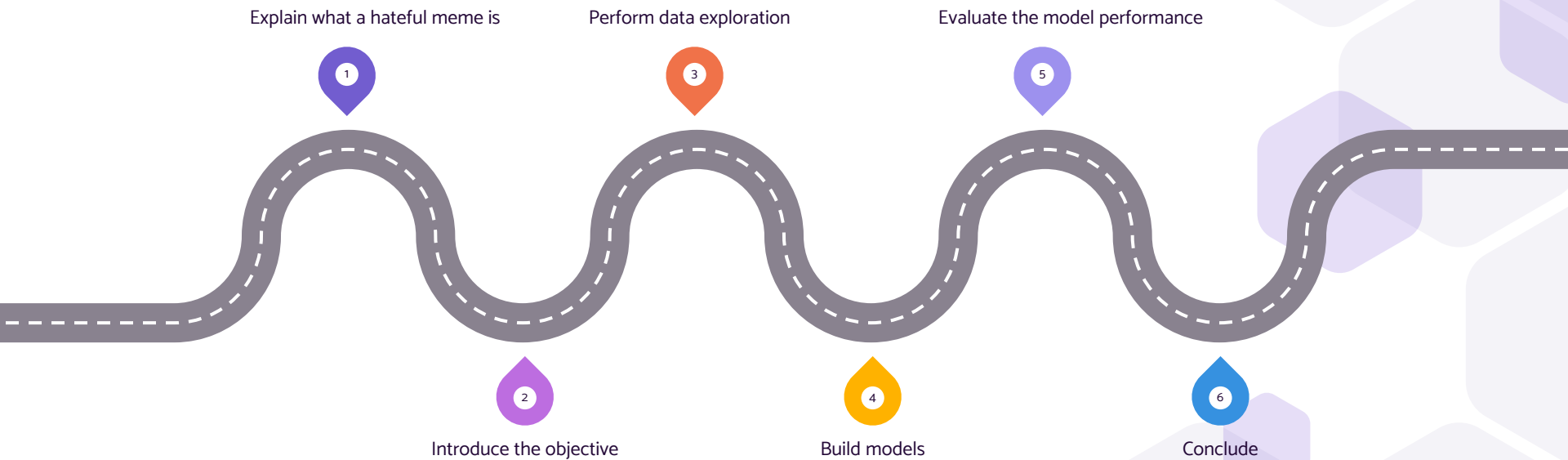
Uyen Huynh

Look how many
people love you

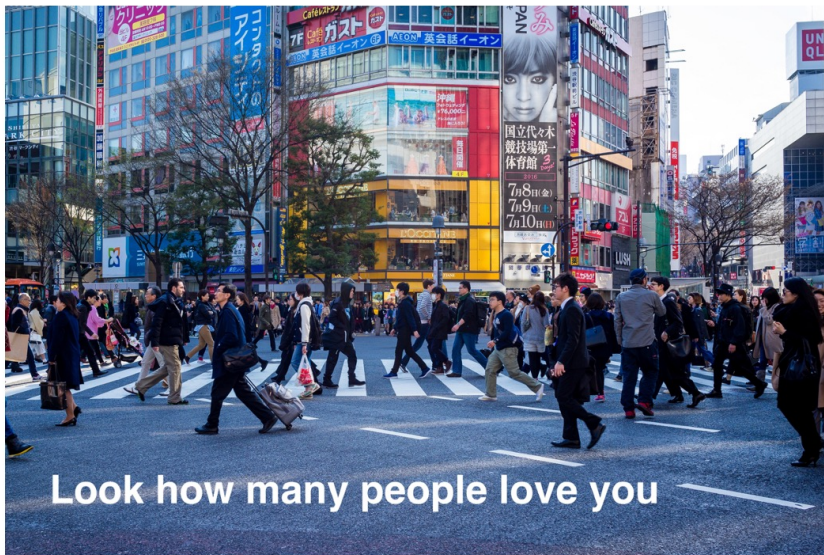




ROADMAP

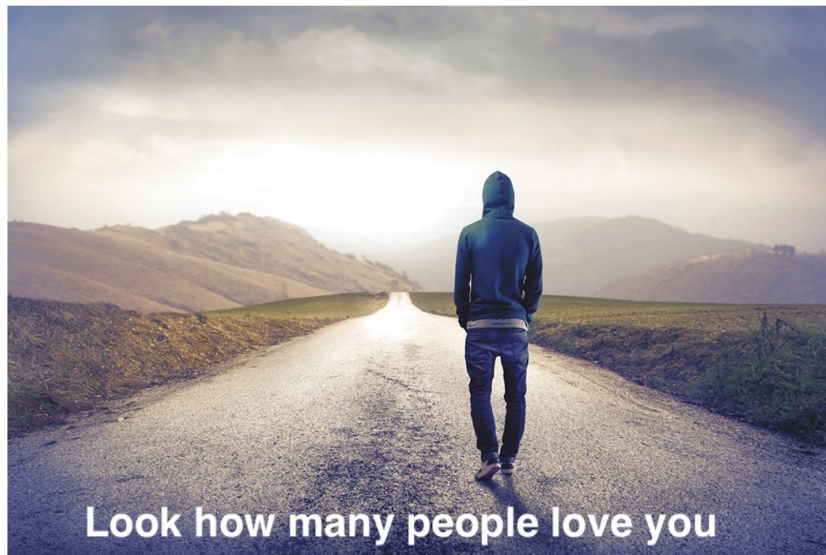


What is a hateful meme?



Look how many people love you

Class 0: Not hateful

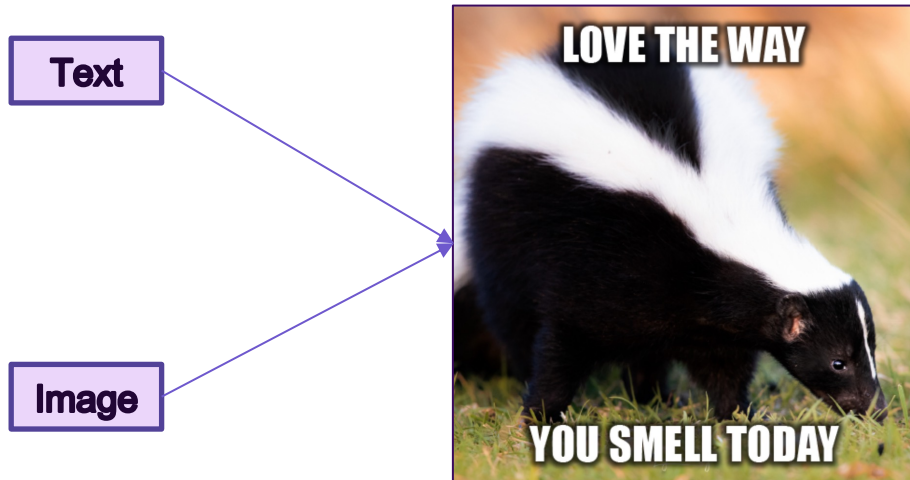


Look how many people love you

Class 1: Hateful



What does a meme comprise of ?

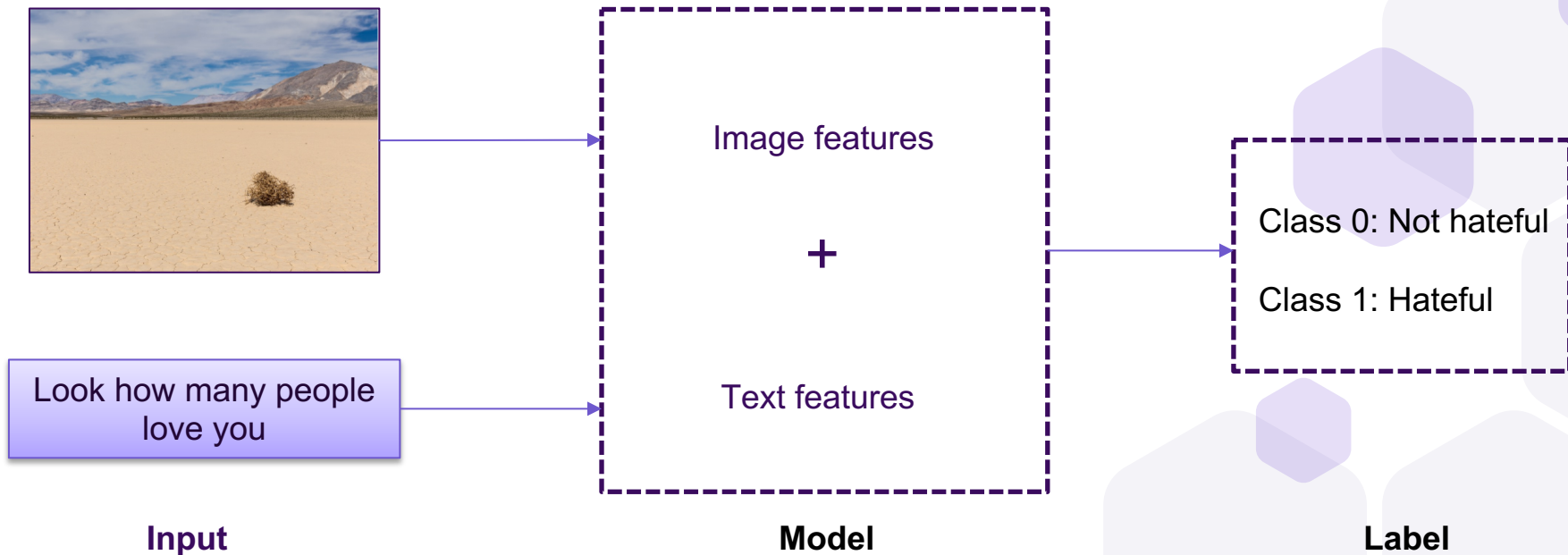


To understand the true meaning of a meme, we can't separately look at the image or text.



Objective

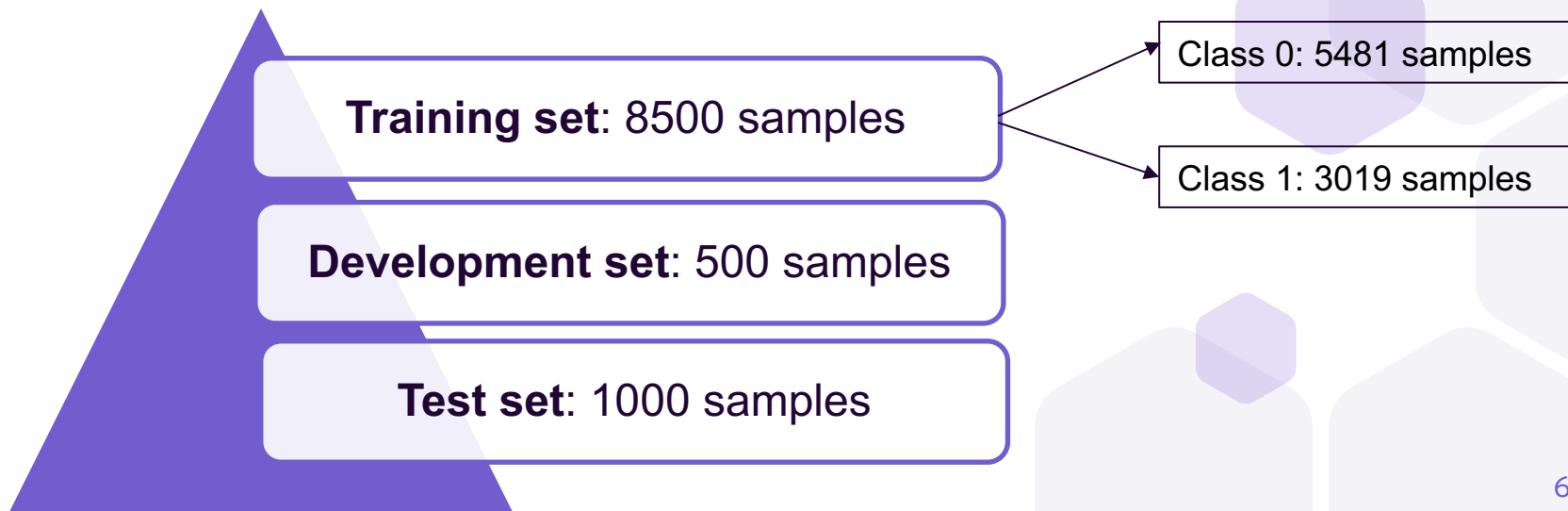
Design a multimodal convolutional neuron network to detect hateful speech in memes



Data

Source: The dataset is created by Facebook AI & can be downloaded on the [Hateful Memes Challenge](#) website.

The dataset contains pairs of text & image:





Data Preparation

Image

- Resize the images to **128 x 128** pixel.
- Normalize the pixel values using the mean and stddev of ImageNet

mean = [0.485, 0.456, 0.406]

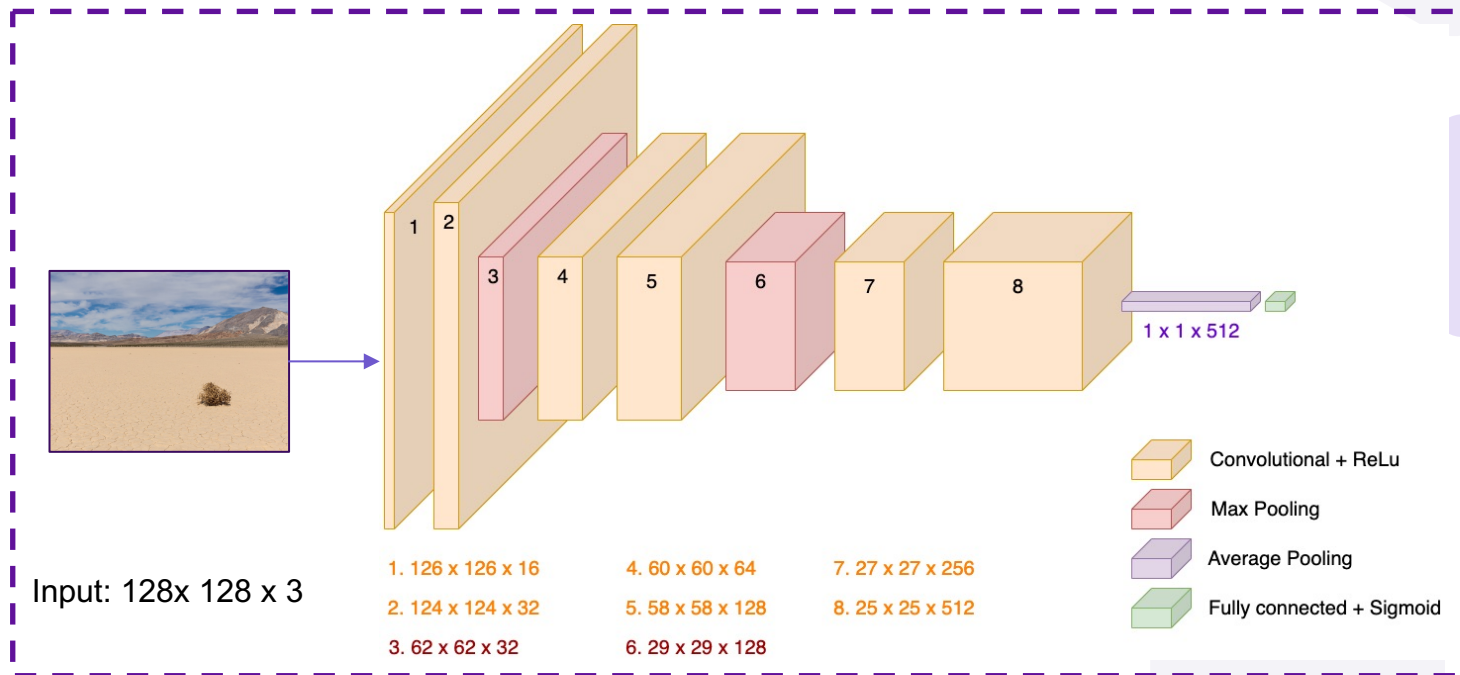
stddev = [0.229, 0.224, 0.225]

Text

- Using NLTK tokenize package to divide the sentence into a list of words.
- Remove stop words.

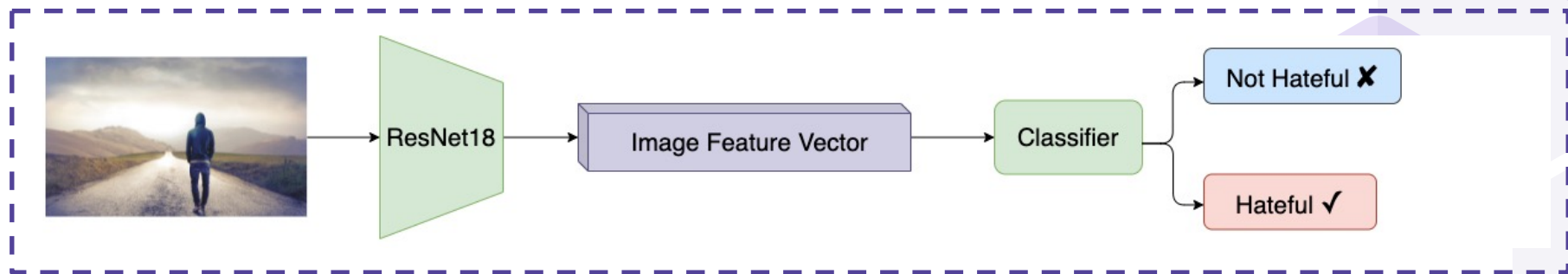
Model: Baseline Unimodal 1

Simple CNN: Using only images to predict the label



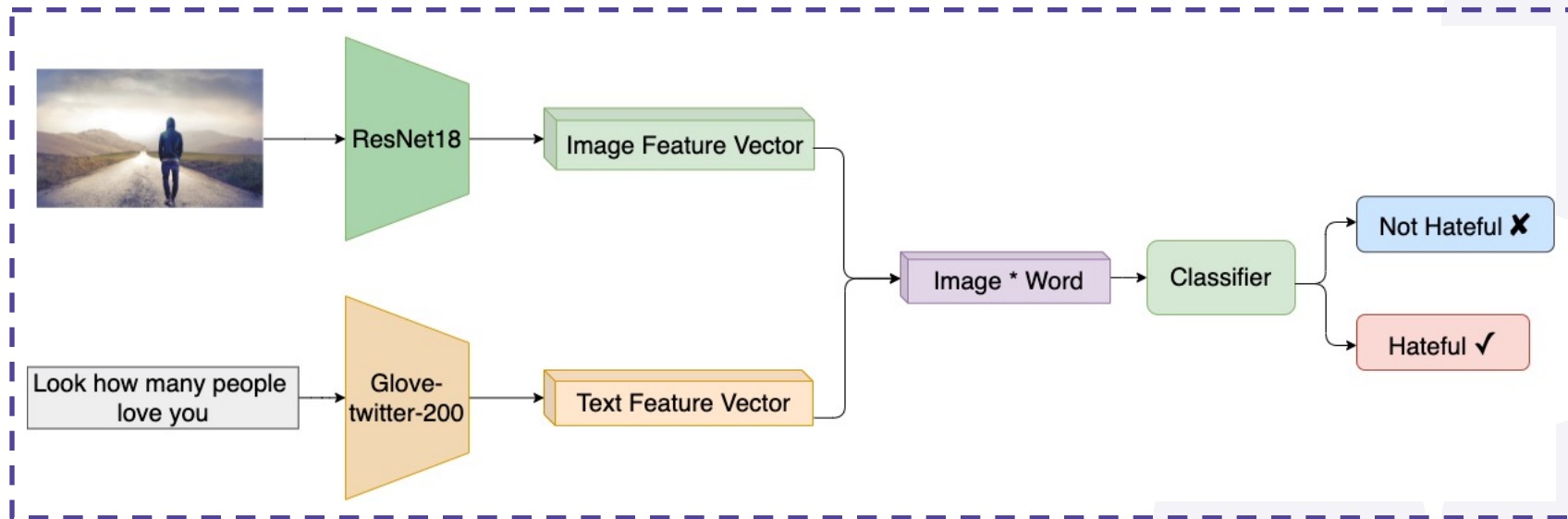
Model: Baseline Unimodal 2

Using pre-trained model ResNet18 to extract image features



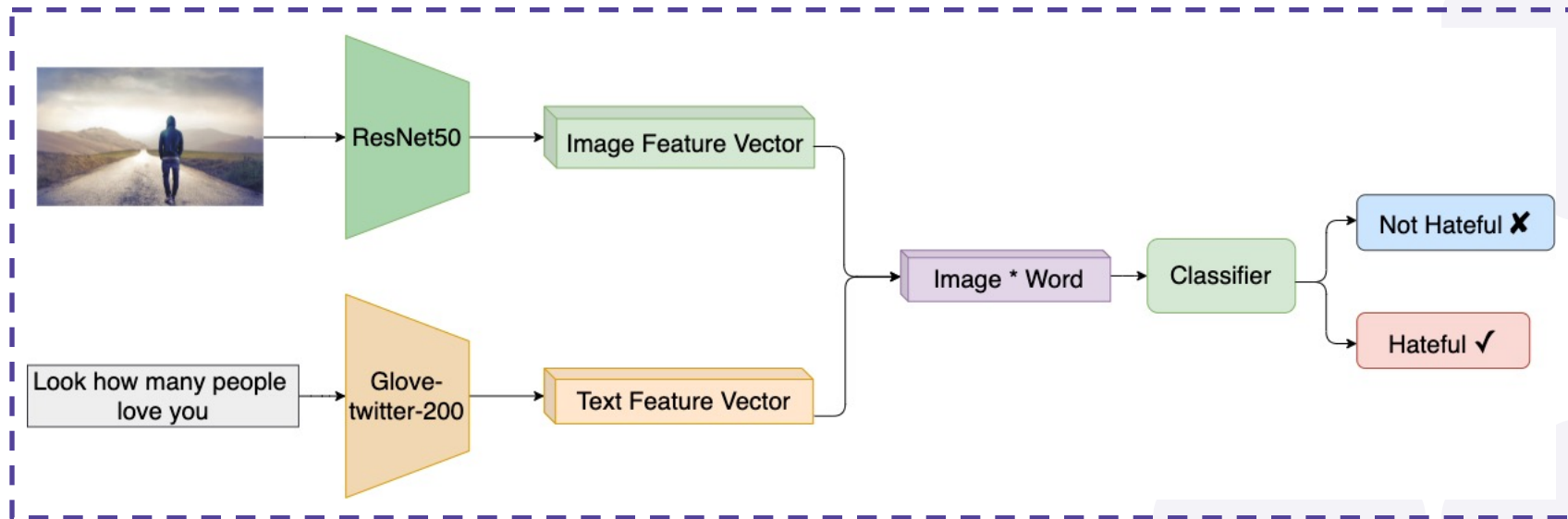
Model: Baseline Multimodal 1

Using pre-trained models: ResNet18 & glove-twitter-200



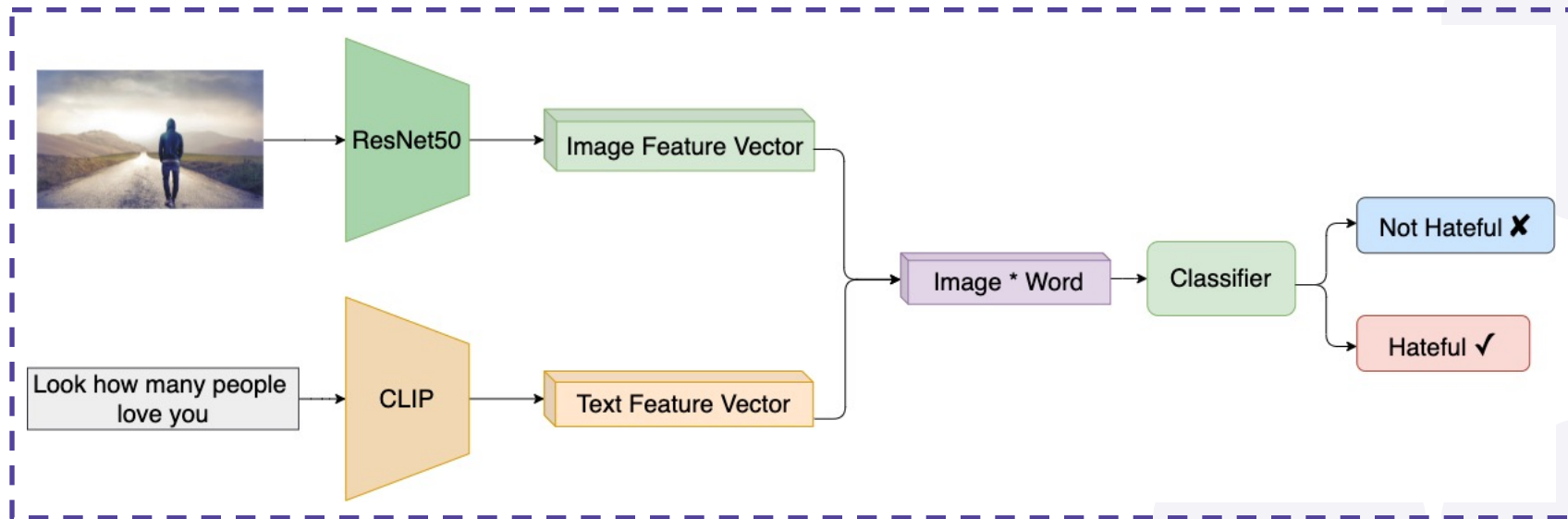
Model: Baseline Multimodal 2

Using pre-trained models: ResNet50 & glove-twitter-200



Model: Text Image Multimodal

Using pre-trained models: ResNet50 & CLIP





Evaluation Metrics

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

TP: True Positive

TN: True Negative

FP: False Positive

FN: False Negative



Model Evaluation

Model	Description	Development Accuracy
Baseline Unimodal 1	Designed CNN	0.506
Baseline Unimodal 2	ResNet18	0.524
Baseline Multimodal 1	ResNet18 & Glove-twitter-200	0.574
Baseline Multimodal 2	ResNet50 & Glove-twitter-200	0.576
Text Image Multimodal	ResNet50 & CLIP	0.616

Text Image 3 has the highest development accuracy

Text Image 3: Optimal Threshold Selection

In the binary classification, the classifier returns a predicted label based on the threshold:

$$\begin{array}{c} \text{Hateful} \\ g(x) \geq \gamma \\ \text{Not Hateful} \end{array}$$

Let's choose the optimal threshold for Text Image 3

```
Threshold = 0.0, Development Accuracy = 0.494
Threshold = 0.1, Development Accuracy = 0.534
Threshold = 0.2, Development Accuracy = 0.570
Threshold = 0.3, Development Accuracy = 0.584
Threshold = 0.4, Development Accuracy = 0.612
Threshold = 0.5, Development Accuracy = 0.616
Threshold = 0.6, Development Accuracy = 0.588
Threshold = 0.7, Development Accuracy = 0.558
Threshold = 0.8, Development Accuracy = 0.526
Threshold = 0.9, Development Accuracy = 0.508
Threshold = 1.0, Development Accuracy = 0.506
Optimal Threshold = 0.5, Development Accuracy = 0.616
```



Text Image 3: Model Performance

Test Accuracy = 0.637

Confusion Matrix

	Predicted		
		Class 0	Class 1
	Actual	Class 0	Class 1
	Class 0	TN = 429	FP = 81
	Class 1	FN = 282	TP = 208

Result Comparison

Type	Model	Validation Accuracy	Test Accuracy
Unimodal	Human		84.70
	Image-Grid	50.67	52.73
	Image-Region	52.53	52.36
	Text BERT	58.27	62.80
Multimodal (Unimodal Pretraining)	Late Fusion	59.39	63.20
	Concat BERT	59.32	61.53
	MMBT-Grid	59.59	62.83
	Text Image	61.60	63.70
	MMBT-Region	64.75	67.66
	ViLBERT	63.16	65.27
	Visual BERT	65.01	66.67

Source: [The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes, FacebookAI](#)



Conclusion

The **Text Image** Model is good at detecting hate speech related to:

- Religion
- Race
- Gender
- Terrorism