



Multimodal Political Bias Identification and Neutralization

By Cedric Bernard, Xavier Pleimling, Amun Kharel and Chase Vickery
CS 6804: Multimodal Vision
Spring 2023
Instructor- Dr. Chris Thomas

Presentation Overview





Introduction



Presentation Overview

- ❖ Introduction and Motivation
- ❖ Background and Related Works
- ❖ Other Debiasing Approaches
- ❖ Proposed Approach
- ❖ Results and Evaluation
- ❖ Future Work, Strengths and Weaknesses

Introduction and Motivation

- ❖ Data from the 2020 election on Growing Political Divide



Introduction and Motivation

- ❖ Data from the 2020 election on Growing Political Divide
- ❖ Pew Research showed that 9 out of 10 people from both parties thought the other candidate getting elected would cause lasting harm



Introduction and Motivation

- ❖ Data from the 2020 election on Growing Political Divide
- ❖ Pew Research showed that 9 out of 10 people from both parties thought the other candidate getting elected would cause lasting harm
- ❖ 8 out of 10 people believe the opposite candidate have diametrically different values





Introduction and Motivation

- ❖ Analysis of 150 million tweets from 3.8 million Twitter Users, pertaining to political and nonpolitical issues



Introduction and Motivation

- ❖ Analysis of 150 million tweets from 3.8 million Twitter Users, pertaining to political and nonpolitical issues
- ❖ Presence of Political Echo Chamber





Textual Bias

❖ Subjective Bias

avan, Mummie, there wouldn't be any railway fare and we shouldn't
oms. Oh, *do* let us go in a caravan."

Mrs. Russell shook her head. "I know it sounds lovely, darling; but
e we to get a caravan? It would cost at least fifty pounds to buy one,
en if we had one, Daddy couldn't get away this summer. No, we
ike up our minds to do without a holiday this year; but I'll tell you what
ll do: we'll all go to Southend for the day, as we did last year, and
r lunch and tea with us and have a splendid picnic."

"Then we can bathe again," said Bob; "but, oh! I do wish I could ha
ny and ride," he added unexpectedly. "You don't know how I long
ny," he continued, sighing deeply as he remembered the blissful holi
en a friend let him share his little Dartmoor pony and ride occasional
"Southend is nothing but houses and people," cried Phyllis; "it's no b
an this place; and oh! Mummie, I do so *long* for fields and flowers
imals," she added piteously; and she shook her long brown hair for
hide the tears in her eyes.

"Never mind, darling, you shall have them one day," answered
assell with easy vagueness.

This really was not very comforting, and it was the most fortunate thing
st at that moment a car stopped at the door.

"Uncle Edward!" shouted Bob, rushing from the room. Phyllis br
e tears so hastily from her eyes that she arrived at the front door almo
on as he did, and both flung themselves on the tall, kindly-looking man st
g beside the car.

"Uncle Edward! Uncle Edward!" they cried. "You've come at
e've been longing to see you. Oh, how glad we are you're here!"

Now the delightful thing was that their uncle seemed just as pleased to
em as they were to see him, and returned their hugs and greetings with
most cordiality. They were just on the point of dragging him into
ouse, hanging one on each arm, when he said: "Stop, not so fast. There
me things to fetch in from the car."

So saying he began diving into the back of it and bringing out, not on
itcase, but various parcels, which he handed out one by one.

"That's the pair of chickens I've brought for your mother," said he, ha



Textual Bias

- ❖ Subjective Bias
- ❖ Emotionally Charged Language (Our focus)

avan, Mummie, there wouldn't be any railway fare and we shouldn't
oms. Oh, *do* let us go in a caravan."

Mrs. Russell shook her head. "I know it sounds lovely, darling; but
e we to get a caravan? It would cost at least fifty pounds to buy one,
en if we had one, Daddy couldn't get away this summer. No, we
ike up our minds to do without a holiday this year; but I'll tell you what
ll do: we'll all go to Southend for the day, as we did last year, and
r lunch and tea with us and have a splendid picnic."

"Then we can bathe again," said Bob; "but, oh! I do wish I could ha
ny and ride," he added unexpectedly. "You don't know how I long
ny," he continued, sighing deeply as he remembered the blissful holi
en a friend let him share his little Dartmoor pony and ride occasional
"Southend is nothing but houses and people," cried Phyllis; "it's no b
an this place; and oh! Mummie, I do so *long* for fields and flowers
imals," she added piteously; and she shook her long brown hair for
hide the tears in her eyes.

"Never mind, darling, you shall have them one day," answered
assell with easy vagueness.

This really was not very comforting, and it was the most fortunate thing
st at that moment a car stopped at the door.

"Uncle Edward!" shouted Bob, rushing from the room. Phyllis bru
e tears so hastily from her eyes that she arrived at the front door almo
on as he did, and both flung themselves on the tall, kindly-looking man st
g beside the car.

"Uncle Edward! Uncle Edward!" they cried. "You've come at
e've been longing to see you. Oh, how glad we are you're here!"

Now the delightful thing was that their uncle seemed just as pleased to
em as they were to see him, and returned their hugs and greetings with
most cordiality. They were just on the point of dragging him into
ouse, hanging one on each arm, when he said: "Stop, not so fast. There
me things to fetch in from the car."

So saying he began diving into the back of it and bringing out, not on
itcase, but various parcels, which he handed out one by one.

"That's the pair of chickens I've brought for your mother," said he, ha



Textual Bias

- ❖ Subjective Bias
- ❖ Emotionally Charged Language (Our focus)
- ❖ Omission of all the facts (Non-trivial problem)

...avan, Mummie, there wouldn't be any railway fare and we shouldn't
oms. Oh, *do* let us go in a caravan."

Mrs. Russell shook her head. "I know it sounds lovely, darling; but
e we to get a caravan? It would cost at least fifty pounds to buy one,
en if we had one, Daddy couldn't get away this summer. No, we
ike up our minds to do without a holiday this year; but I'll tell you what
ll do: we'll all go to Southend for the day, as we did last year, and
r lunch and tea with us and have a splendid picnic."

"Then we can bathe again," said Bob; "but, oh! I do wish I could ha
ny and ride," he added unexpectedly. "You don't know how I long
ny," he continued, sighing deeply as he remembered the blissful holi
en a friend let him share his little Dartmoor pony and ride occasional
"Southend is nothing but houses and people," cried Phyllis; "it's no b
an this place; and oh! Mummie, I do so *long* for fields and flowers
imals," she added piteously; and she shook her long brown hair for
hide the tears in her eyes.

"Never mind, darling, you shall have them one day," answered
Russell with easy vagueness.

This really was not very comforting, and it was the most fortunate thing
st at that moment a car stopped at the door.

"Uncle Edward!" shouted Bob, rushing from the room. Phyllis bru
e tears so hastily from her eyes that she arrived at the front door almo
on as he did, and both flung themselves on the tall, kindly-looking man st
g beside the car.

"Uncle Edward! Uncle Edward!" they cried. "You've come at
e've been longing to see you. Oh, how glad we are you're here!"

Now the delightful thing was that their uncle seemed just as pleased to
em as they were to see him, and returned their hugs and greetings with
most cordiality. They were just on the point of dragging him into
ouse, hanging one on each arm, when he said: "Stop, not so fast. There
me things to fetch in from the car."

So saying he began diving into the back of it and bringing out, not on
itcase, but various parcels, which he handed out one by one.

"That's the pair of chickens I've brought for your mother," said he, ha

Image Bias

- ❖ Facial Expression and correlation with Bias



Image Bias

- ❖ Facial Expression and correlation with Bias
- ❖ Replacing biased images with neutralized images on the same topic





Related Works



Automatically Neutralizing Subjective Bias in Text

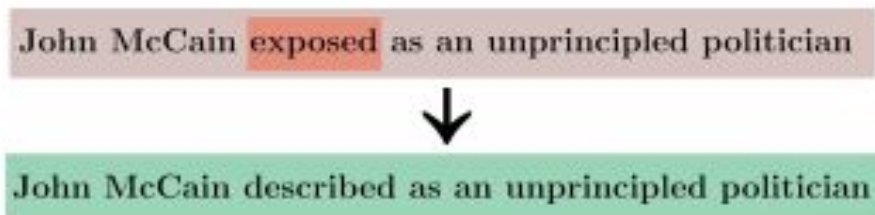
John McCain exposed as an unprincipled politician



John McCain described as an unprincipled politician



Automatically Neutralizing Subjective Bias in Text



Goal: detect and edit biased text in a controllable manner



Automatically Neutralizing Subjective Bias in Text

Proposed Solutions:

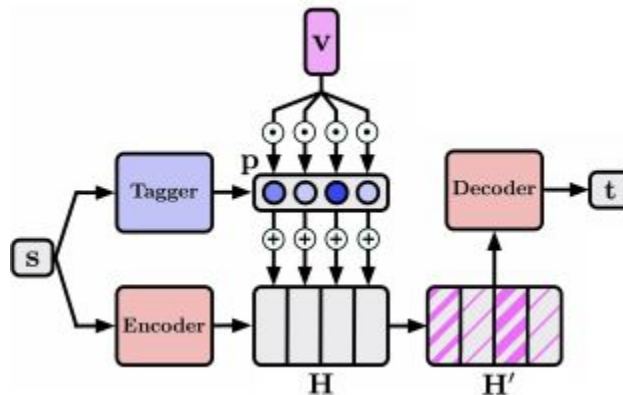
1. MODULAR
2. CONCURRENT

Automatically Neutralizing Subjective Bias in Text

Proposed Solutions:

1. MODULAR
2. CONCURRENT

MODULAR

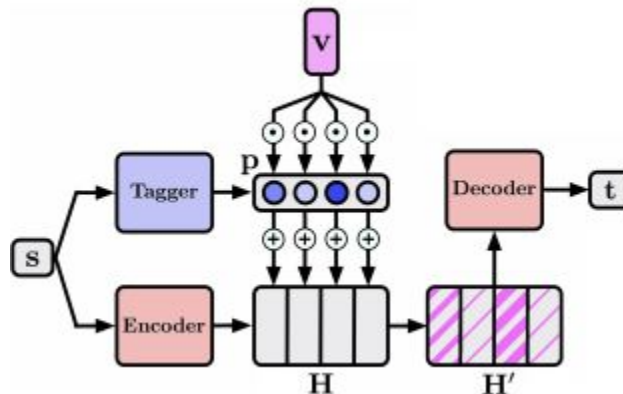


Automatically Neutralizing Subjective Bias in Text

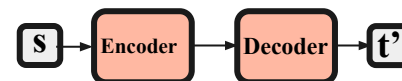
Proposed Solutions:

1. MODULAR
2. CONCURRENT

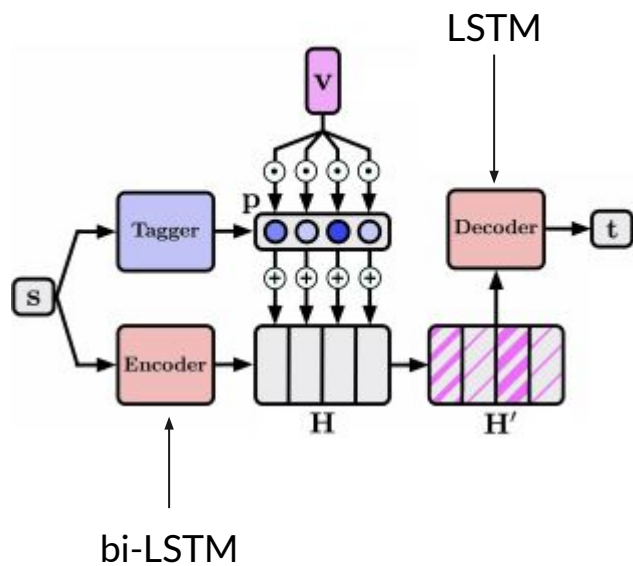
MODULAR



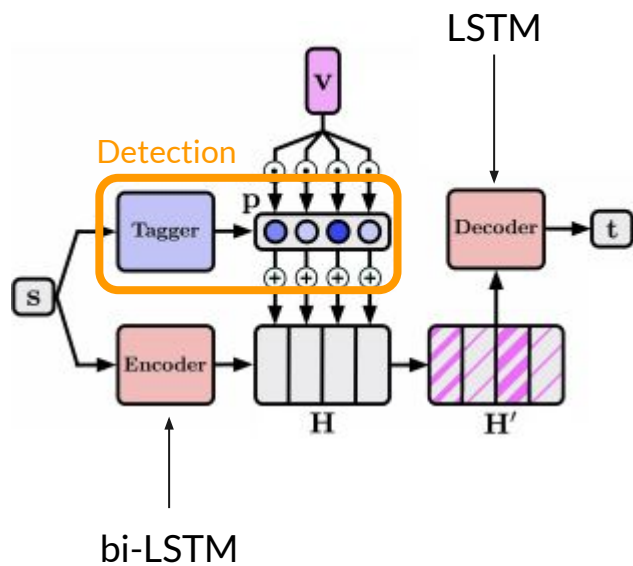
CONCURRENT



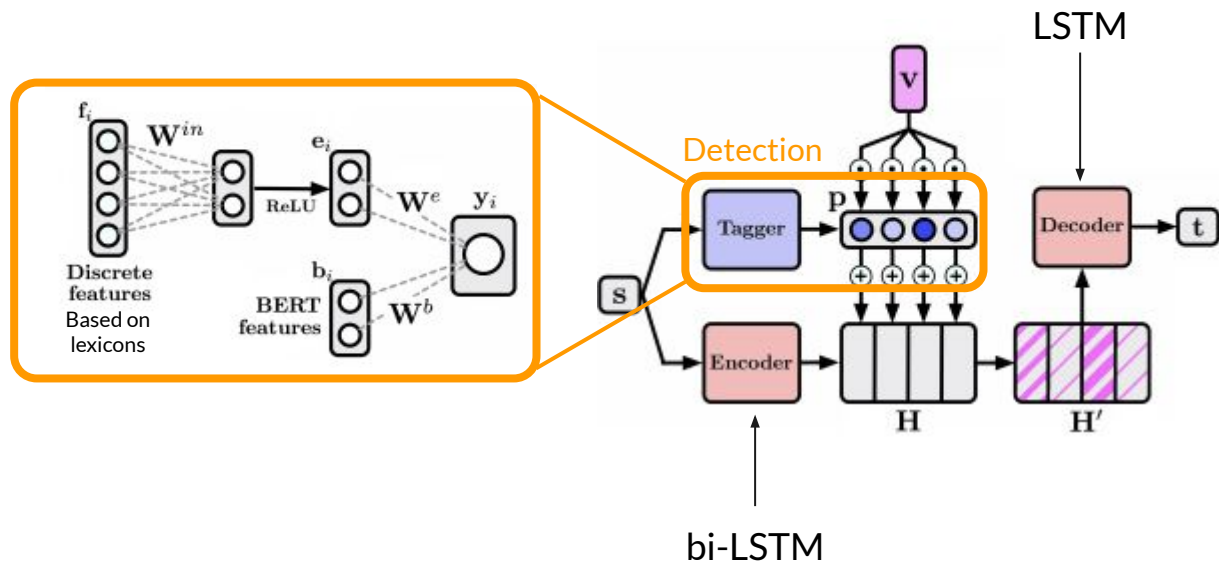
MODULAR



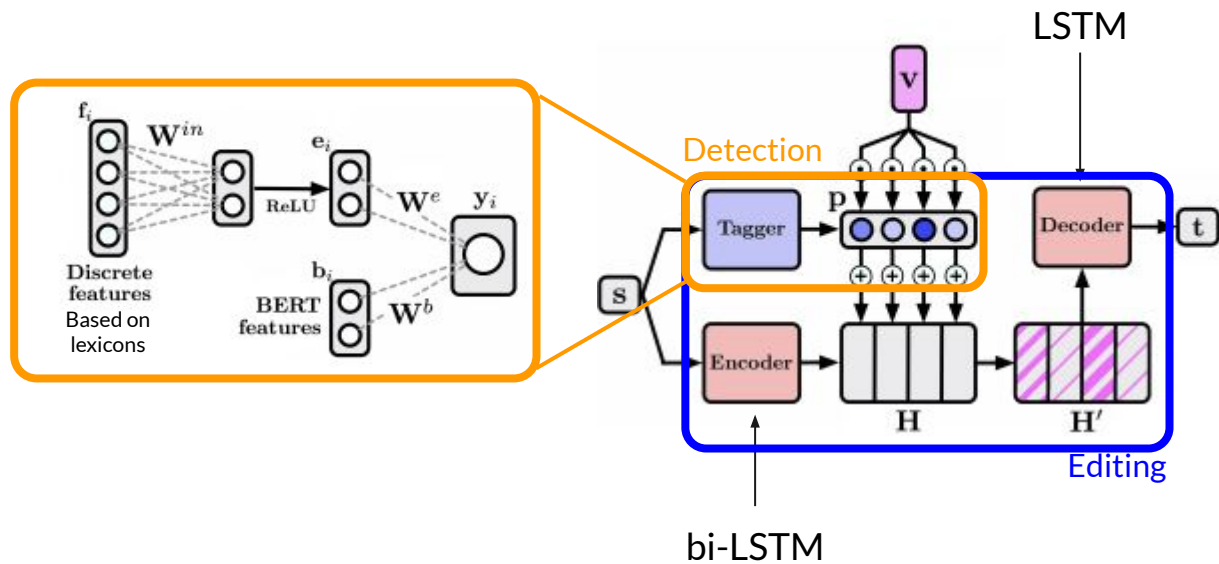
MODULAR



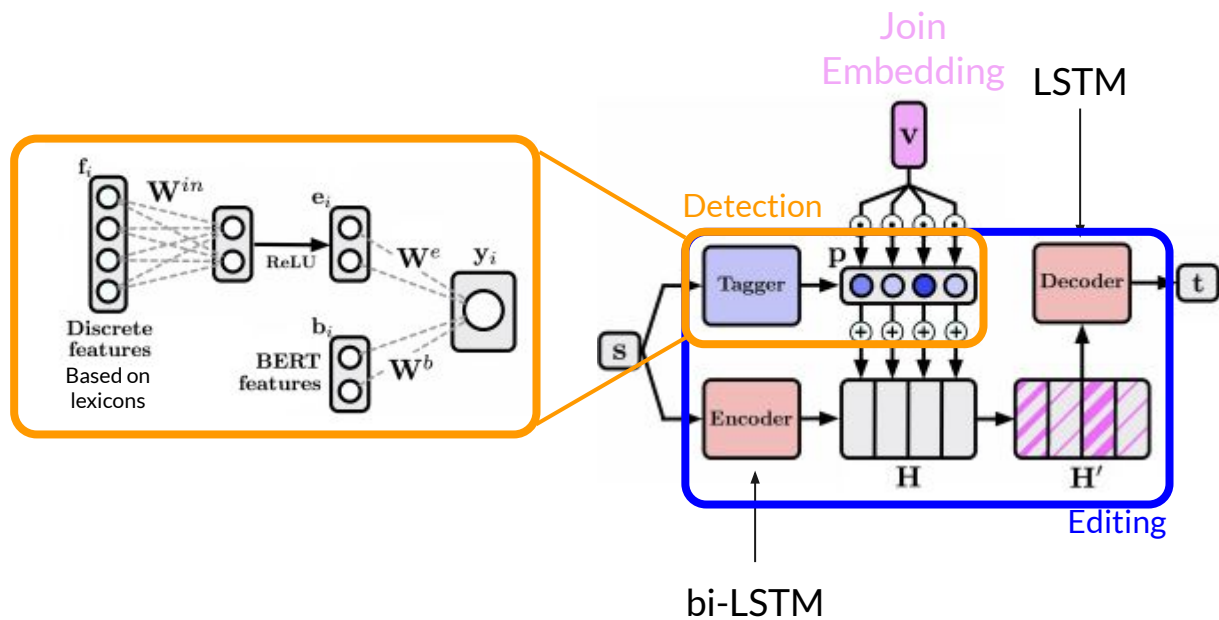
MODULAR



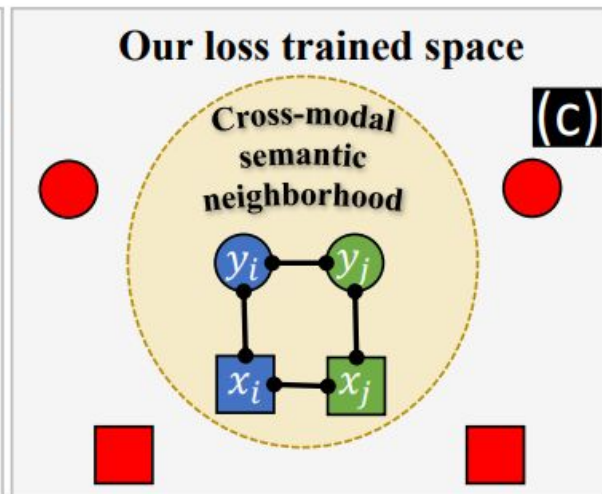
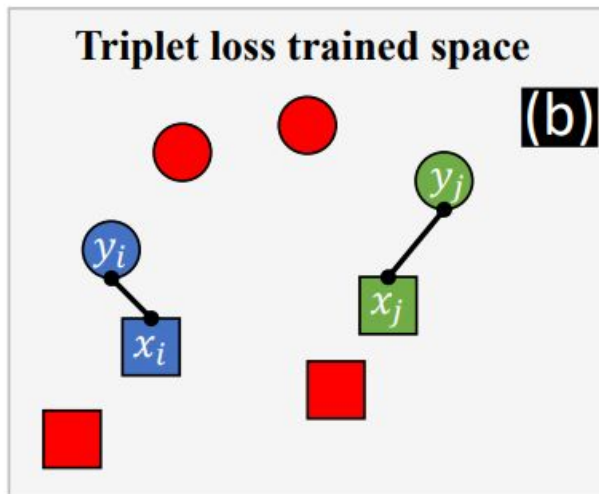
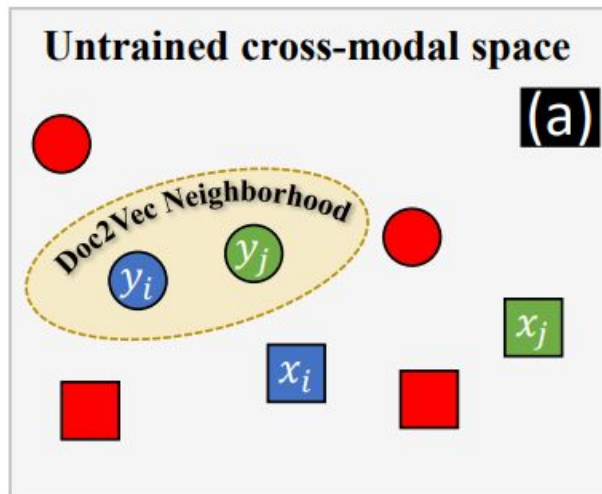
MODULAR



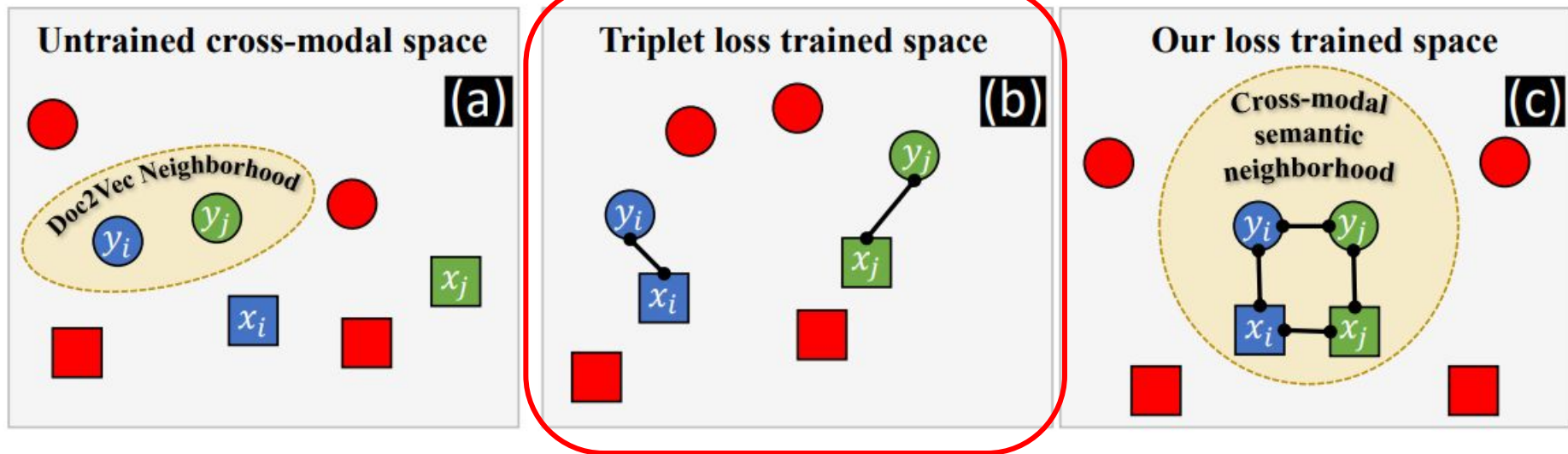
MODULAR



Semantic Neighborhoods

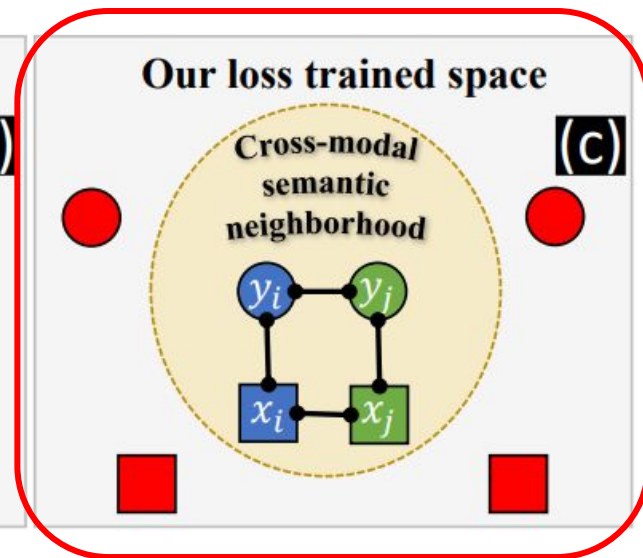
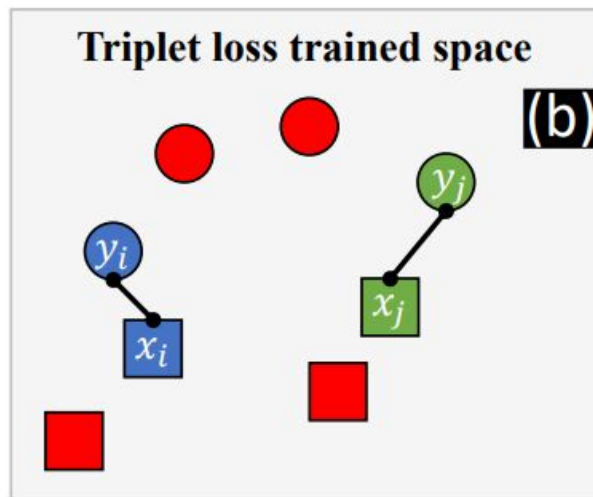
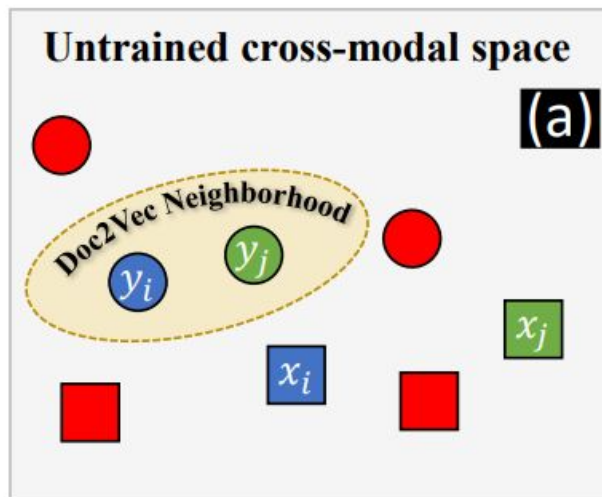


Semantic Neighborhoods



$$\mathcal{L}_{ang}^{NP+SYM}(\mathcal{T}, \mathcal{T}_{sym}) = \mathcal{L}_{ang}^{NP}(\mathcal{T}) + \mathcal{L}_{ang}^{NP}(\mathcal{T}_{sym})$$

Semantic Neighborhoods



$$\mathcal{L}_{text}(\mathcal{T}'_{text}) = \mathcal{L}_{ang}(\mathcal{T}'_{text})$$

$$\mathcal{L}_{img}(\mathcal{T}'_{img}) = \mathcal{L}_{ang}(\mathcal{T}'_{img})$$



Bias Identification in Text

Several techniques used:

- Pre-trained BERT/models for detection and editing of bias
- Adversarial training
- Attention-based mechanisms for detection
- Gaussian mixture models



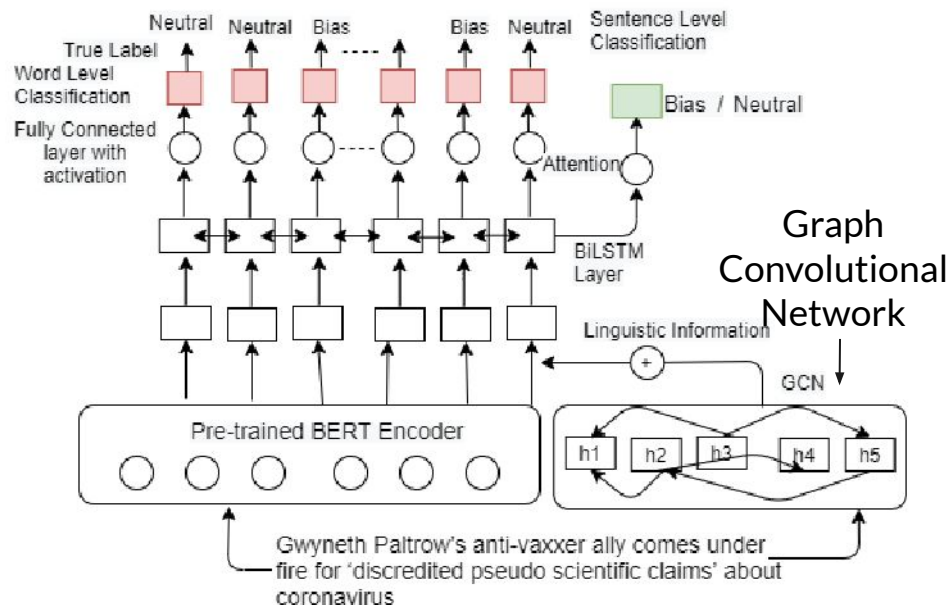
Bias Identification in Text

Several techniques used:

- Pre-trained BERT/models for detection and editing of bias
- Adversarial training
- Attention-based mechanisms for detection
- Gaussian mixture models

These approaches only focus on text data and thus do not take advantage of leveraging visual data

Other Debiasing Approaches



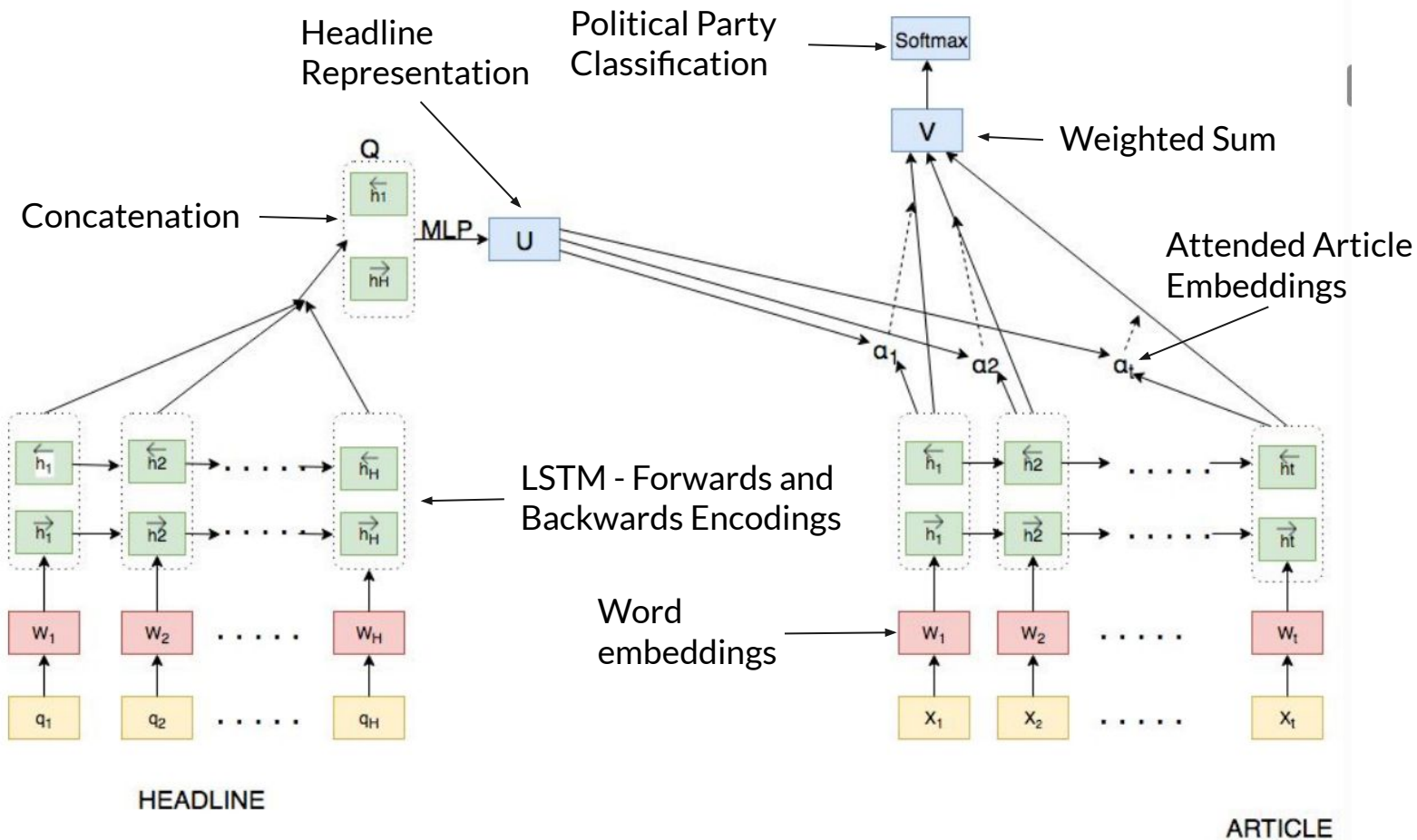
$$\min_D \frac{1}{N} \sum_{i=1}^N \mathcal{L}(D(H(x_i)), z_i) \quad (2)$$

H = Encoder

$$\min_{H,C} \frac{1}{N} \sum_{i=1}^N \alpha \cdot \mathcal{L}(C(H(x_i)), y_i) + (1 - \alpha) \cdot \mathcal{L}(D(H(x_i)), 0.5) \quad (3)$$

C = Classifier

D = Decoder





Bias Score Annotation



Politics Dataset

- Article Text - Image Pairs
- Left vs. Right
- 20 Topics



Labelling Bias Score

- ❖ Each website given score of -1 to 1
- ❖ -1 being Far Left, 0 being neutral and 1 being Far right
- ❖ For example: Score of -0.5 will be considered moderately left leaning

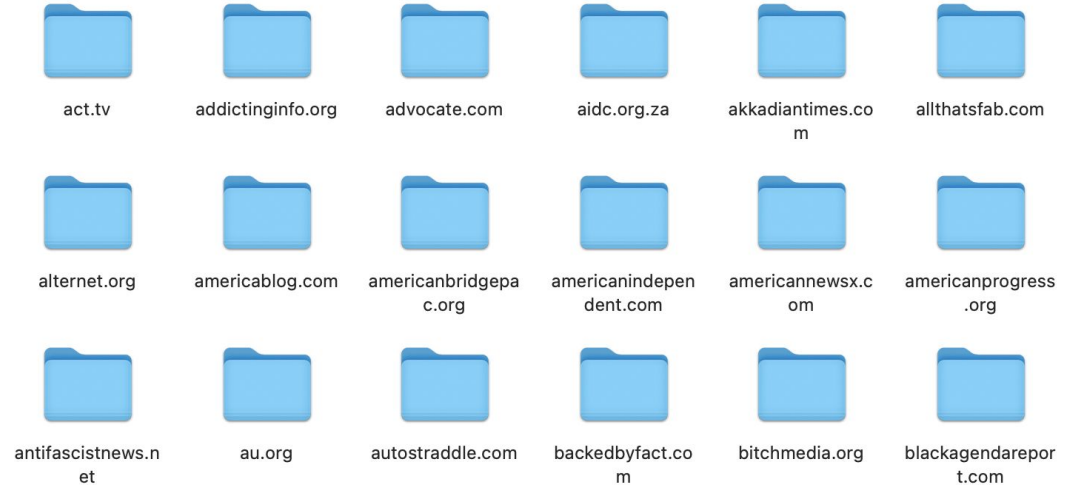
Labelling Bias Score

- ❖ We have 20 different topics in the Politics Dataset images



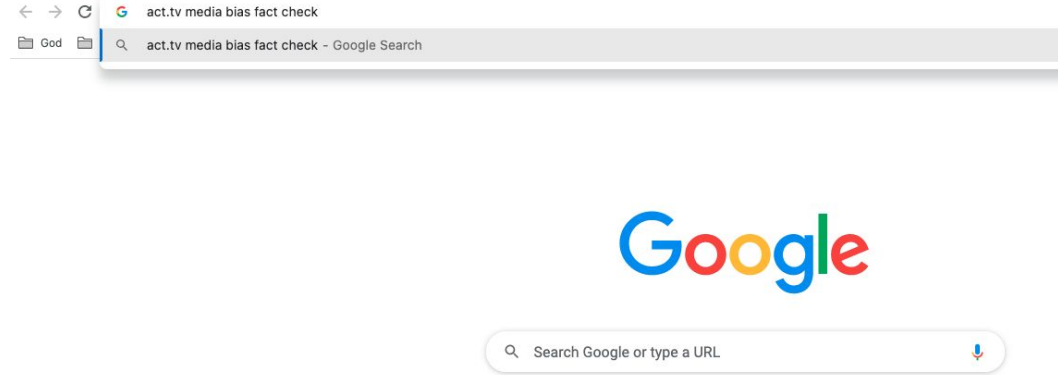
Labelling Bias Score

- ❖ We have the same websites under each topic for left leaning websites and right leaning website



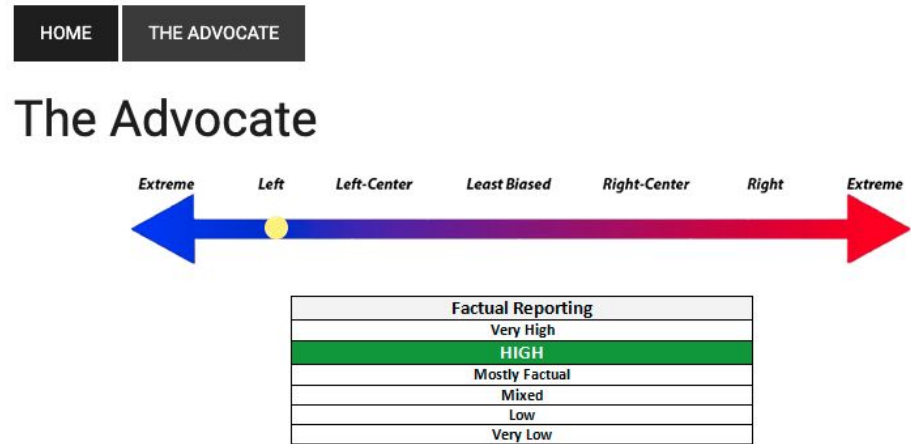
Labelling Bias Score (Media Bias Fact Check)

- ❖ Google the website with Media Bias Fact Check added to it



Labelling Bias Score

- ❖ From mediabiasfactcheck.com we retrieve the score.



Labelling Bias Score

- ❖ Get the score and put it on spreadsheet

	A	B	C
1	Name	Leaning	Bias
2	act.tv	left	-0.8
3	addictinginfo.org	left	-0.9
4	advocate.com	left	-0.95
5	aidc.org.za	left	-0.3
6	akkadiantimes.com	left	-0.5
7	allthatsfab.com	left	-0.85
8	alternet.org	left	-0.9
9	americablog.com	left	-0.5
10	americanbridgepac.org	left	-0.6
11	americanindependent.com	left	-0.5
12	americannewsx.com	left	-0.7
13	americanprogress.org	left	-0.5
14	antifascistnews.net	left	-0.3
15	au.org	left	-0.3
16	autostraddle.com	left	-0.5
17	backedbyfact.com	left	-0.25
18	bitchmedia.org	left	-0.45
19	blackagendaareport.com	left	-0.95
20	blacklivesmatter.com	left	-0.45
21	blue-route.org	left	-0.6
22	bluedotdaily.com	left	-1
23	bluenationreview.com	left	-0.9
24	boingboing.net	left	-0.5



Labelling Bias Score (Exceptions)

- ❖ 0-2 = Least Biased
- ❖ 2-5 = Left/Right Center Bias
- ❖ 5-8 = Left/Right Bias
- ❖ 8-10 = Extreme Bias



Labelling Bias Score (Exceptions)

The categories are as follows:

- 1) Biased Wording/Headlines- Does the source use loaded words to convey emotion to sway the reader. Do headlines match the story?
- 2) Factual/Sourcing- Does the source report factually and back up claims with well-sourced evidence?
- 3) Story choices: Does the source report news from both sides, or do they only publish one side?
- 4) Political Affiliation: How strongly does the source endorse a particular political ideology? Who do the owners support or donate to?



Labelling Bias Score (Exceptions)

Here is an example of how CNN scored and why they were placed in the middle of Left Bias:

Biased Wording= 4 (CNN uses moderate biased words that favor liberals, and headlines typically match the story)

Factual/Sourcing=5 (CNN has failed fact checks and sometimes omits critical information from stories to favor their perspective)

Political Affiliation = 5 (CNN's ownership owns other left-leaning outlets and favors Democratic Candidates)



Labelling Bias Score (Exceptions)

Total = 23

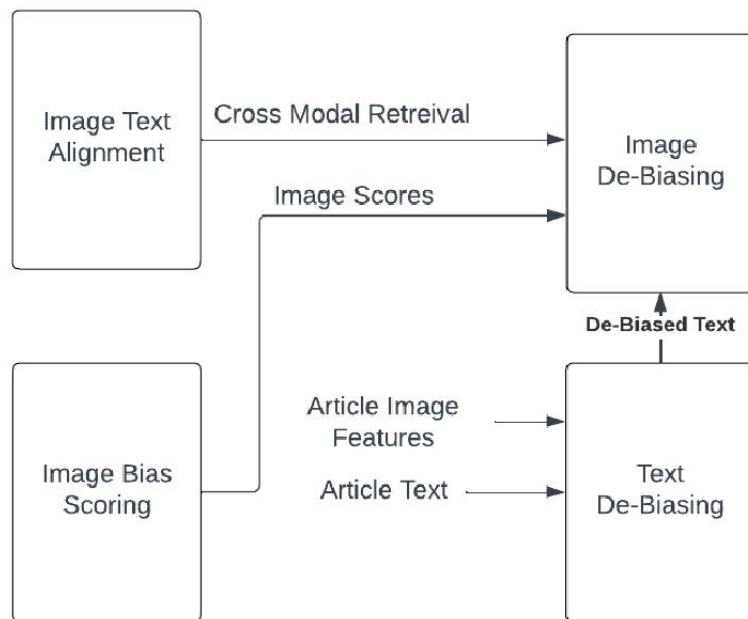
Average = $23/4 = 5.75$

5.75 = Moderate Left Bias

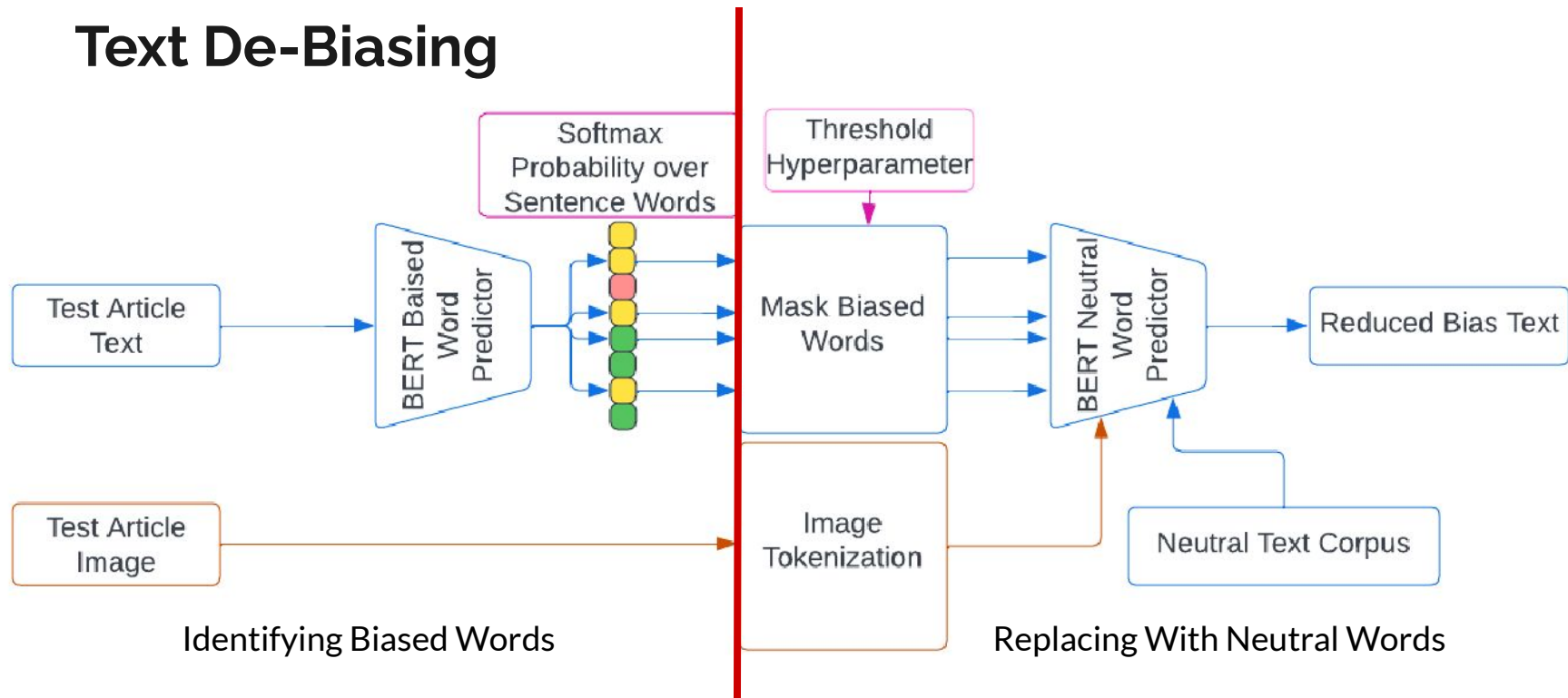


Proposed Approach

Proposed Approach Overview

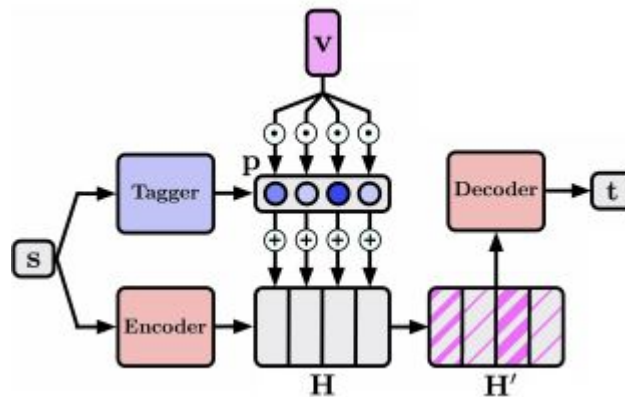


Text De-Biasing



Biased Word Identification

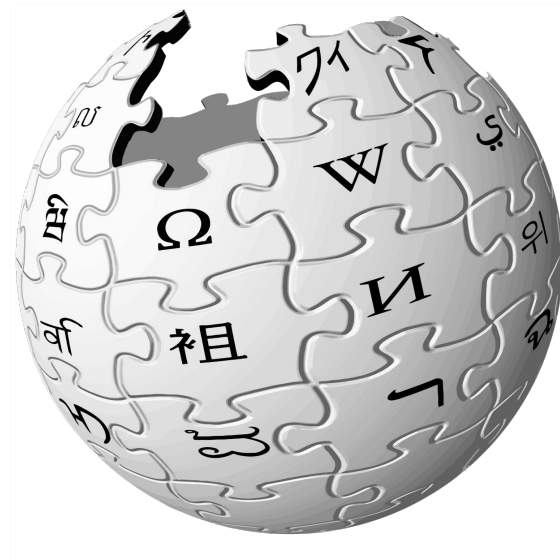
- ❖ Adopting the MODULAR detection module
- ❖ Trained with the Wikipedia Neutrality Corpus (WNC)
- ❖ Returns a set of probabilities that are transferred to the Editing portion of the model



Wikipedia Image Text Dataset

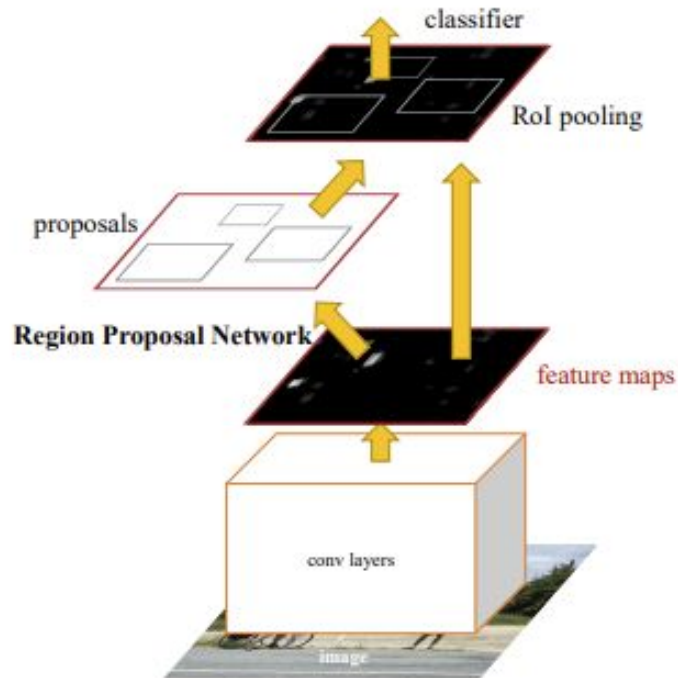
- ❖ 3.5 Million paired Image and text samples
 - Using 1% (35,000)
- ❖ Scraped from text and images from Wikipedia
- ❖ Images compressed to 124,124,3
- ❖ Used as neutral data inputs

[wit/DATA.md at main · google-research-datasets/wit \(github.com\)](#)



[Wikipedia](#)

Faster RCNN



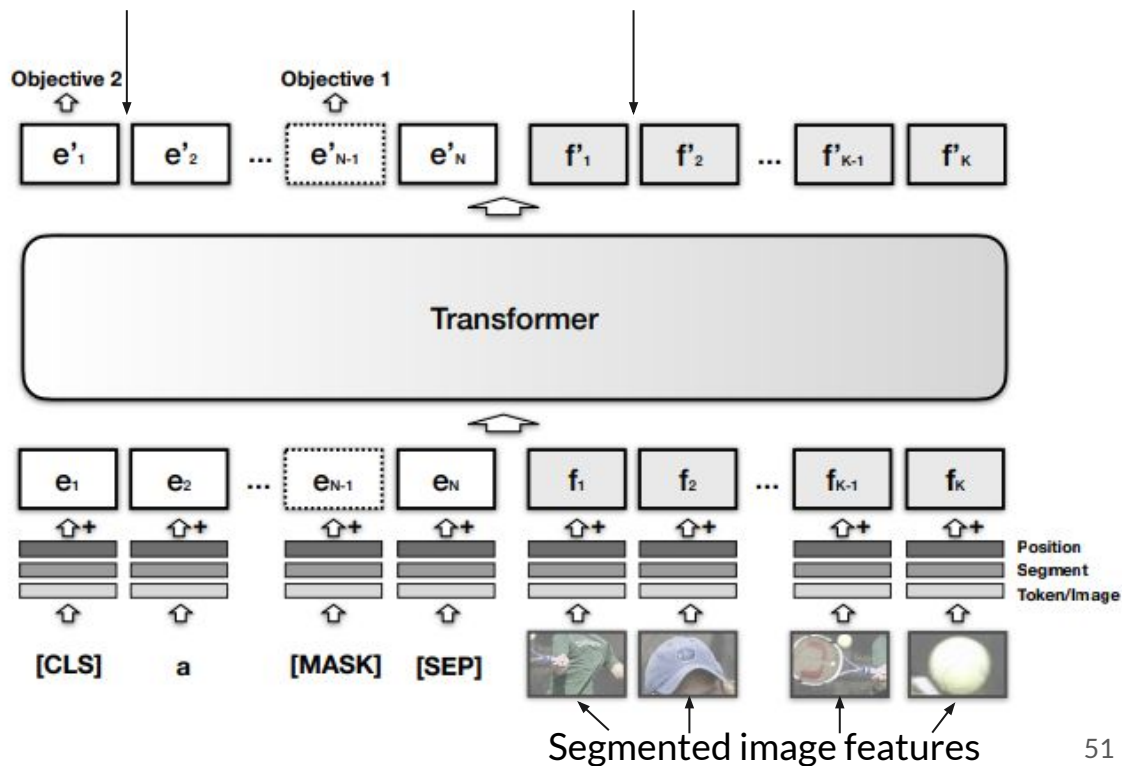
Visual BERT



A person hits a ball with a tennis racket

Sentence-Image Prediction

Masked-Language Modeling



Train Time Word Prediction

Human migration is
the movement of
people from one place
to another



15% Random
Mask

CNN Feature
Extraction

VisualBERT

[Human migration - Wikipedia](#)

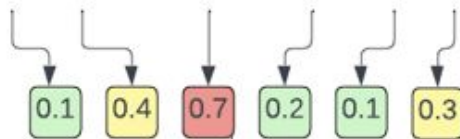


Text Predictor Model Parameters

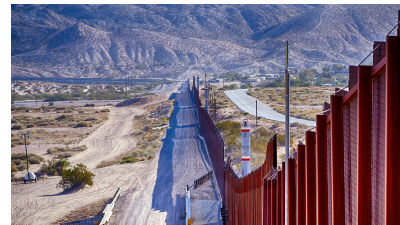
- Faster RCNN ~ 65M parameters (frozen)
- VisualBERT ~ 112M parameters, 3704 trainable parameters
 - Freeze all but last four layers which contains the weights and bias for the predictor and sequential relationship

Test Time Word Prediction

Three migrants caught crossing the border



Three migrants [MASK] crossing the border



CNN Feature
Extraction

Image De-Biasing

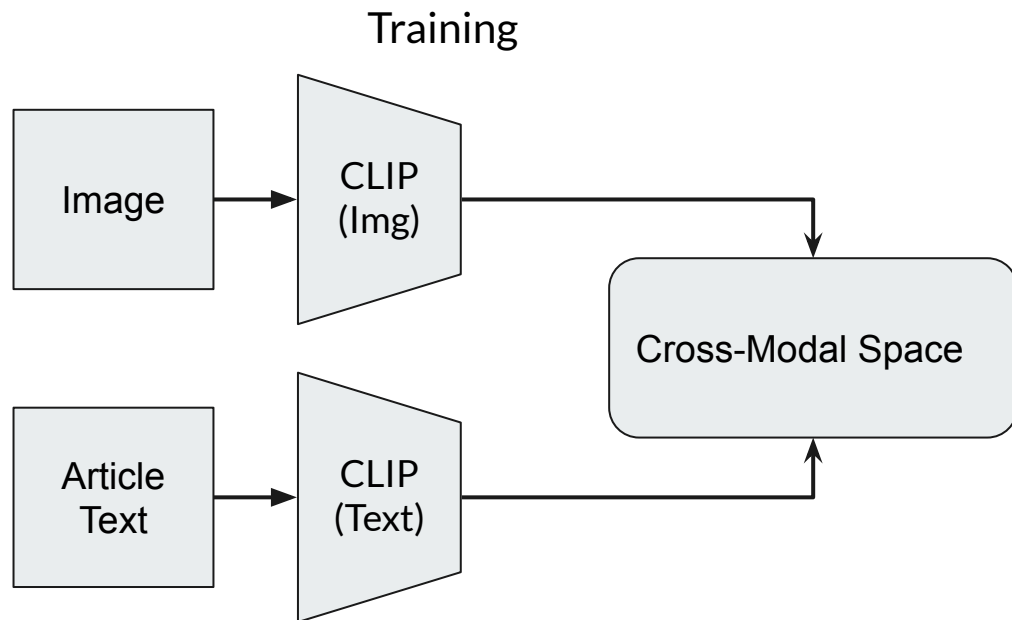


Image De-Biasing

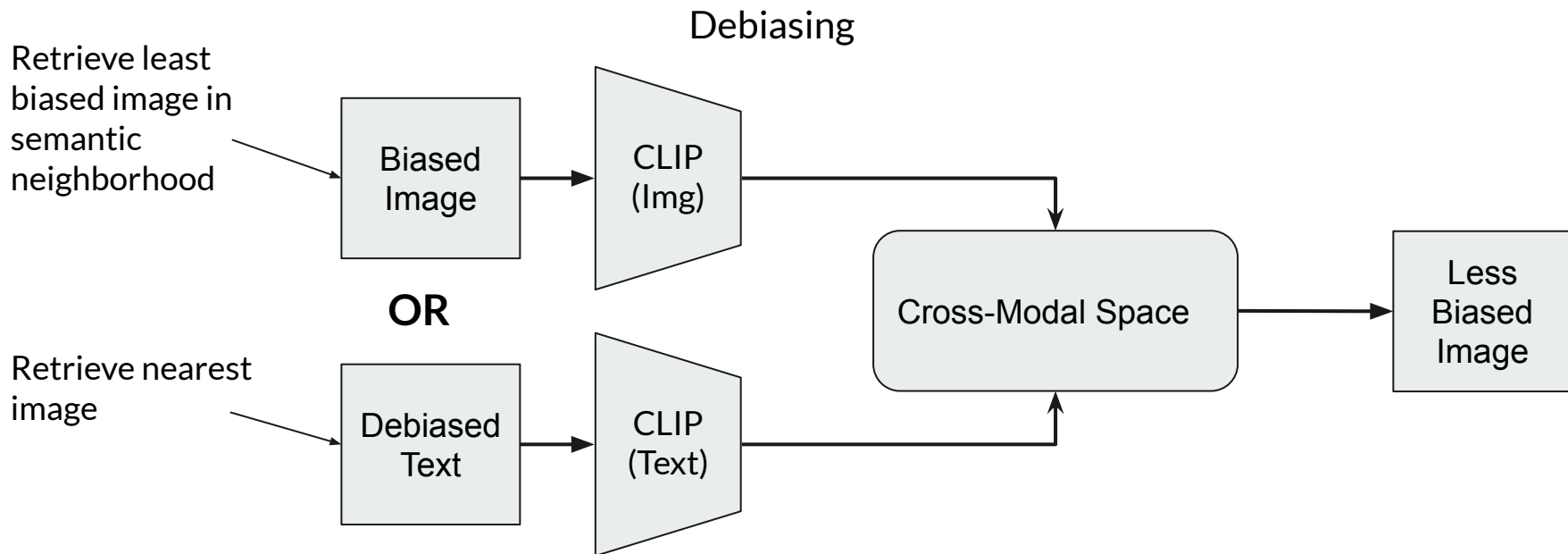


Image De-Biasing

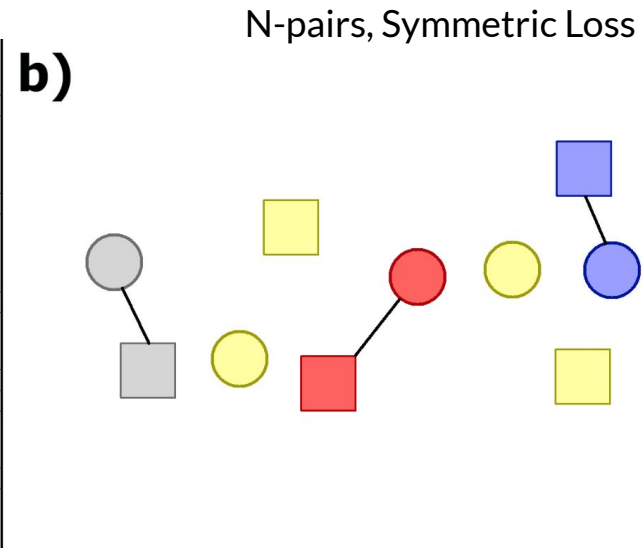
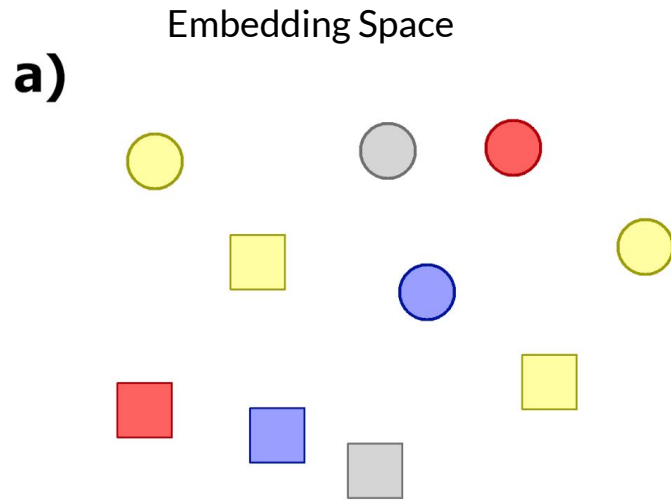


Image De-Biasing

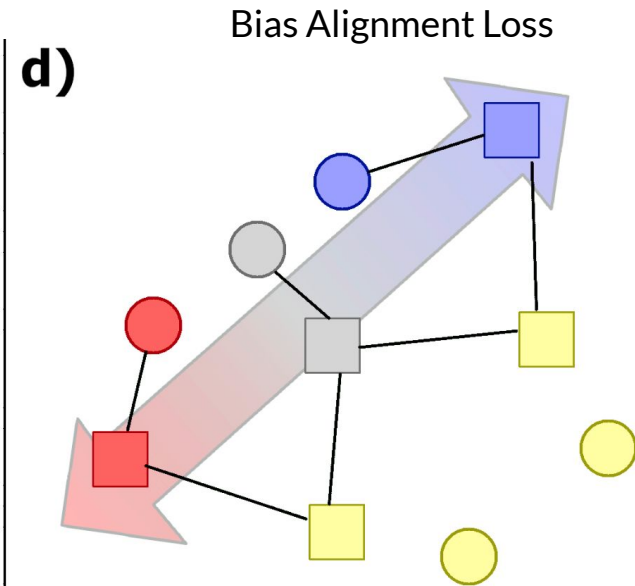
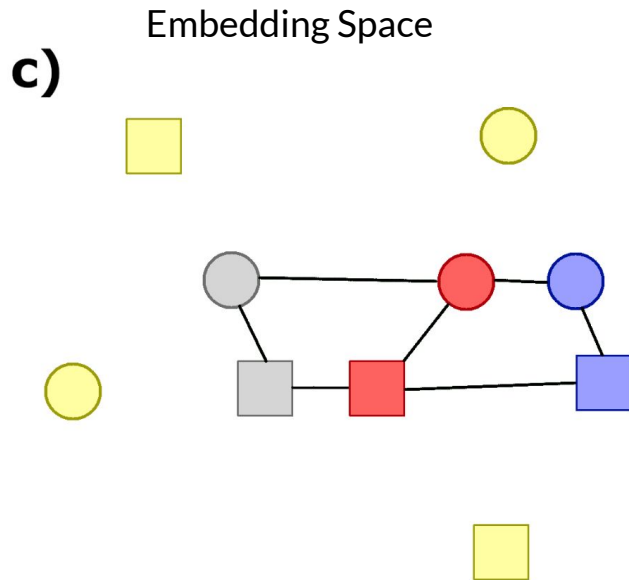
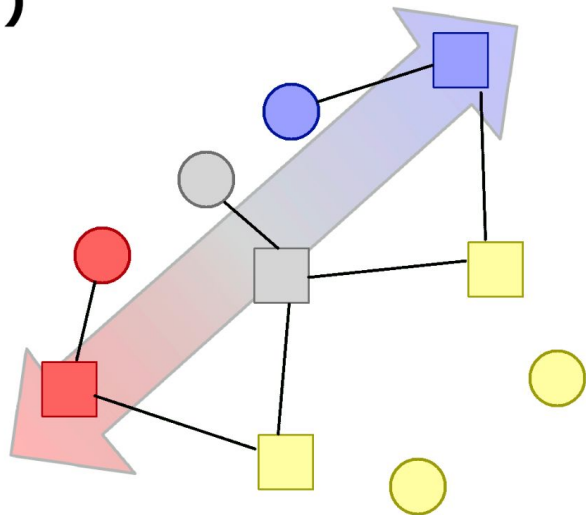


Image De-Biasing

d)



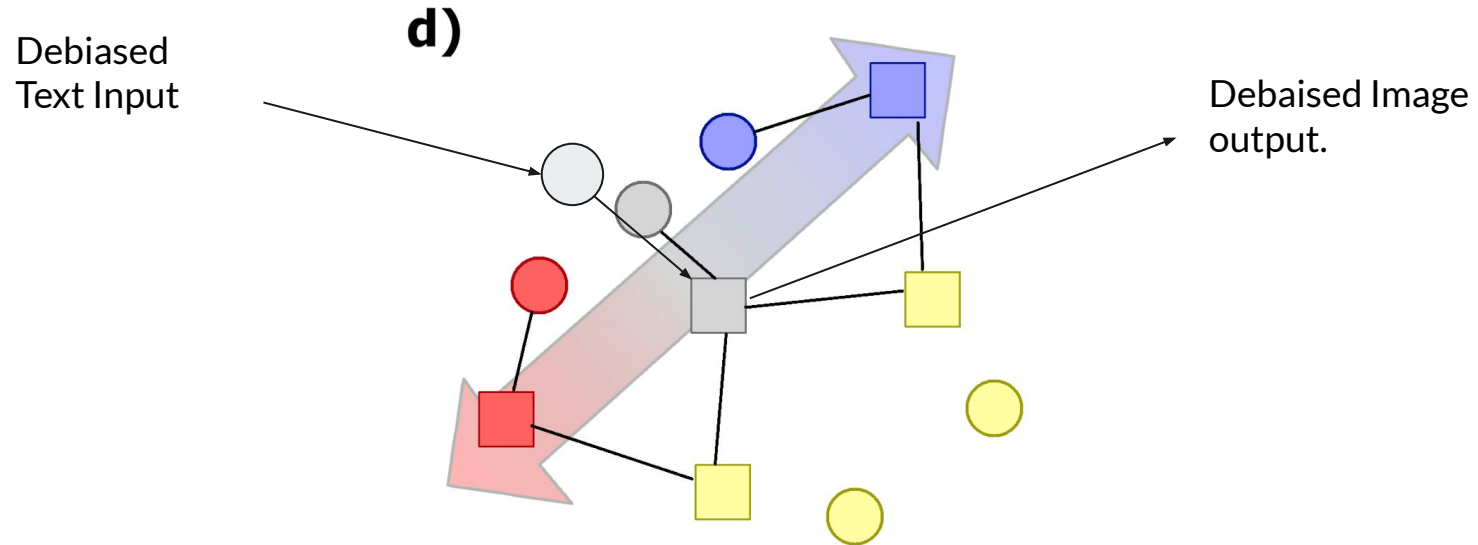
Anchor Bias Neighborhood Negative Sample

$$\mathcal{T}'_{img} = (x_i^a, x_{i'}^p, x_j^n)$$

$$\mathcal{L}_{img}(\mathcal{T}'_{img}) = \mathcal{L}_{ang}(\mathcal{T}'_{img})$$

$$\mathcal{L}_{ang}(\mathcal{T}'_{img}) = \left[\|x_i^a - x_{i'}^p\|_2^2 - 4 \tan^2 \alpha \|x_j^n - \mathcal{C}_i\|_2^2 \right]_+$$

Image De-Biasing (Retrieval)

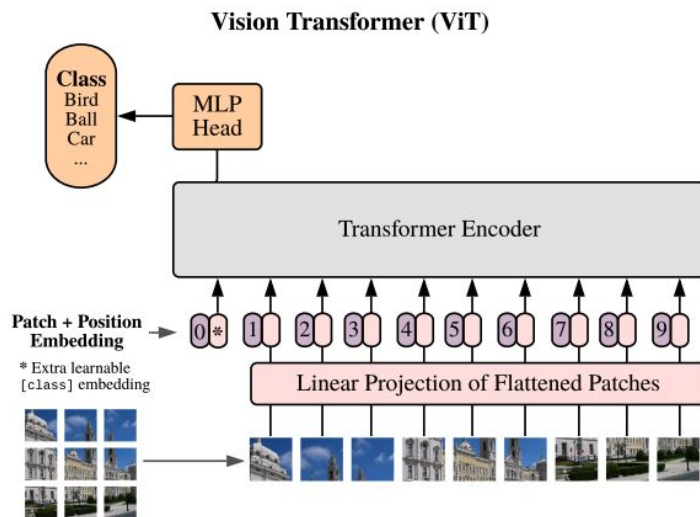




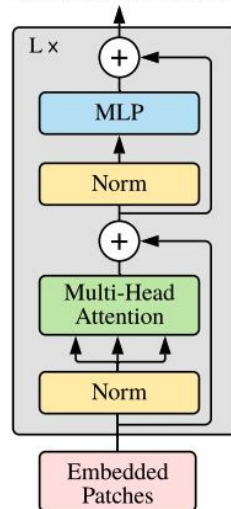
ViT Alignment (Overview)

- ❖ Divide an input image into 196 (14 X 14) small images of size (16 X 16)
- ❖ Treat it as embedding in NLP
- ❖ Use it as input for traditional transformer encoder (like in BERT)
- ❖ We modified it to 2 instead of 12 Transformer layers because of the size of our data and resources we had available
- ❖ Take the last output, use it as input for Dense layer with one linear output (original paper 1000 classes)

ViT Alignment (Overview)



Transformer Encoder





Why Don't We Use A Full Image for Transformer?

- ❖ Because of n^2 parameters
- ❖ Equates to 22 billion parameters only for 224 X 224 X 3 image



Our modifications to the model

- ❖ Our Dataset has 224 X 224 images with respective bias score
- ❖ Changed MLP head to output a linear classifier instead of multi-class classifier
- ❖ We use Mean Squared Loss instead of Cross Entropy Loss
- ❖ We transformed from 12 layers to 2 Transformer Layers
- ❖ We changed MLP size to 512



Results and Evaluation



Biased Word Identification - Results

Qualitative Analysis

Biased Word Identification - Results

Qualitative Analysis

Main example from paper

Input

John McCain exposed as an unprincipled politician

Tokens

(descending order of bias)

exposed	un
politician	##pled
##ci	as
##pr	john
##in	an
mccain	

Biased Word Identification - Results

Qualitative Analysis

Randomly selected from WNC test set

Input

ali once dated american actor jonathan brandis
who unfortunately died in 2003.

Tokens

(descending order of bias)

.	dated	died
2003	once	in
american	actor	
ali	brand	
unfortunately	who	
jonathan	##is	

Biased Word Identification - Results

Input

The leader of the centre left SPD party, currently playing a supporting role as the junior coalition partner in the German government, hopes the capitalise on voter weariness to prevent Mrs Merkel gaining a fourth term in office. He is due to announce whether he will officially run at the end of this month, but Bild announced on Tuesday that he is the candidate.

Tokens

(descending order of bias)

fourth	office	whether	the
the	partner	german	of
##ise	announced	the	a
capital	tuesday	due	mrs
term	bi	will	as
that	coalition	voter	role
is	##ld	officially	government
.	he	this	...

Biased Word Identification - Results

Input

The Wall around Jerusalem (Photo: Musa Al-Shaer)
GENEVA -- The Commission on Human Rights today
started its general debate on the violation of human rights
in the occupied Arab territories, including Palestine, after it
concluded its discussion on the right to development.
Many speakers denounced Israel's actions in the occupied
Palestinian territory, in particular the construction of a
separation Wall.

Tokens

(descending order of bias)

right	s	-	the
particular	territory	the	arab
rights	the	,	palestine
rights	the	the	.
development	occupied	human	violation
discussion	.	its	-
its	occupied	human	:
the	general	actions	...



Biased Word Replacement - Results

- Cosine similarity between predicted masked words and original word embeddings. Word embeddings obtained through pre-trained Gensim
- Human evaluation for bias measurement



Image Debias Evaluation

$$\frac{1}{|Y|} \sum_{y \in Y} |b(N(y))|$$

Bias

Nearest image of input.

Text test set

$$\frac{1}{|Y|} \sum_{x, y \in X, Y} (|b(x)| - |b(N(y))|)$$

Image test set



ViT Model evaluation (RMSE)

- ❖ Stands for Root Mean Square Error (RMSE)
- ❖ Standard Deviation of the residuals (prediction errors).

$$\text{RMSE}_{fo} = \left[\sum_{i=1}^N (z_{fi} - z_{oi})^2 / N \right]^{1/2}$$

Where:

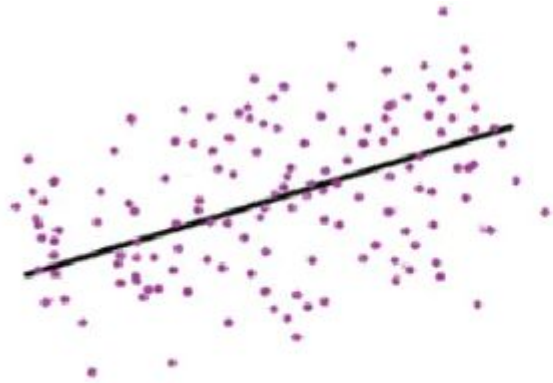
- Σ = **summation** (“add up”)
- $(z_{fi} - z_{oi})^2$ = differences, squared
- N = **sample size**.



ViT Model evaluation (R-squared Score)

- ❖ Regression Score function
- ❖ Negative score- Model can be arbitrarily worse
- ❖ Zero score- Model that does not explain the variability of the response data around its mean
- ❖ Best possible is 1.0, which corresponds to a model that explains the variability of the the response variable around its mean

ViT Model evaluation (R-squared Score)



R-squared : 17%



R-squared : 83%



Human Evaluation

- Form to evaluate 4 key questions comparing novel image and text with original
 - If novel image and text are semantically similar to original
 - If novel text has reduced bias
 - If novel image has reduced bias
 - If novel image and novel text are still aligned

https://docs.google.com/forms/d/e/1FAIpQLSfG3aR9B0VTDJvGKLPN4vy8mqRkBoFz59l5VmXgbSaXAYWoWQ/viewform?usp=sf_link



Strengths and Weaknesses

- Weaknesses
 - Reduced bias text may be less truthful
 - Stacked models may compound error
 - Overall model may be sensitive to rare words, punctuation, etc.
 - People watch audio and video sources and our method does not have neutralizing video and audio.
- Strengths
 - Considers both text and visual information for bias reduction
 - Image debias is a retrieval task. Should make debiasing fast.
 - Our paper, if successful in its implementation, will give people an opportunity to consider viewing news that are neutralized. Also, since we are bombarded by news through the internet, we can analyze the effects of neutralized news and its effect on a person's well being.



Ethical Considerations

- When would debiasing an article's images or text be considered censorship?
- If the negative or positive aspects being covered are truthful, is it appropriate to replace an image that reflects that emotional element with one that does not?
- Public perception of model itself.



Future Work

- Examine “directionality” of bias in embedding space. New loss likely makes regions of neutral representations. Generation?
- Weight bias loss on difference b/w biases (Higher bias difference leads to greater separation).
- Image debiasing metric: Do they still match the context & semantics?
- Image debiasing metric: Measure variety of retrieved images.

Future Work (Timestamp Videos)

