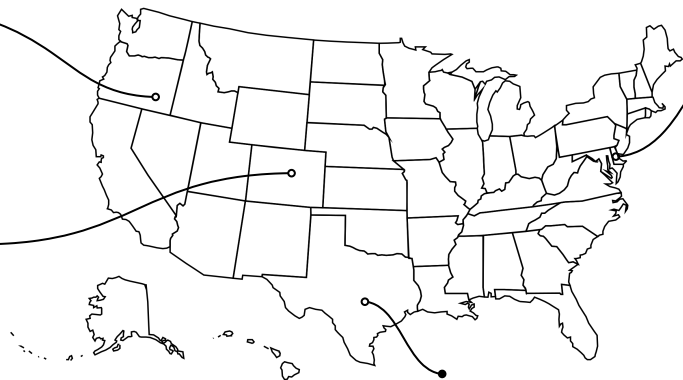


Multimodal Modeling for Political Science Research



@xxxx: F*** Trump. #FamiliesBelongTogether

@xxxx:
#FamiliesBelongTogether
in Denver this morning.



@fams2gether: Nobody,
no matter who you are,
where you are from, you
should spend an hour like
this!



@ResistSnow: #FreeTheChildren
#FamiliesBelongTogether
Our children long for realistic maps of
the future that they can be...



Andreu Casas

Royal Holloway Univ. of London

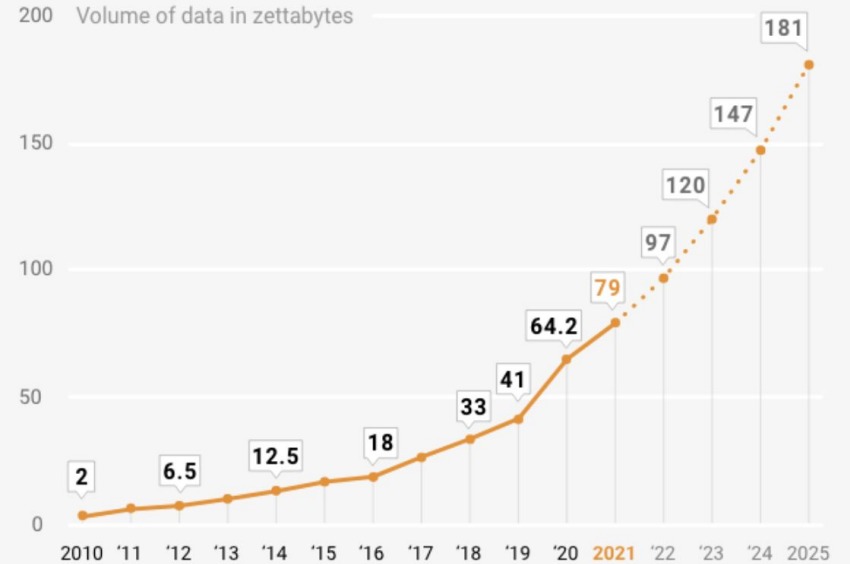
Freek Cool

Vrije Universiteit Amsterdam

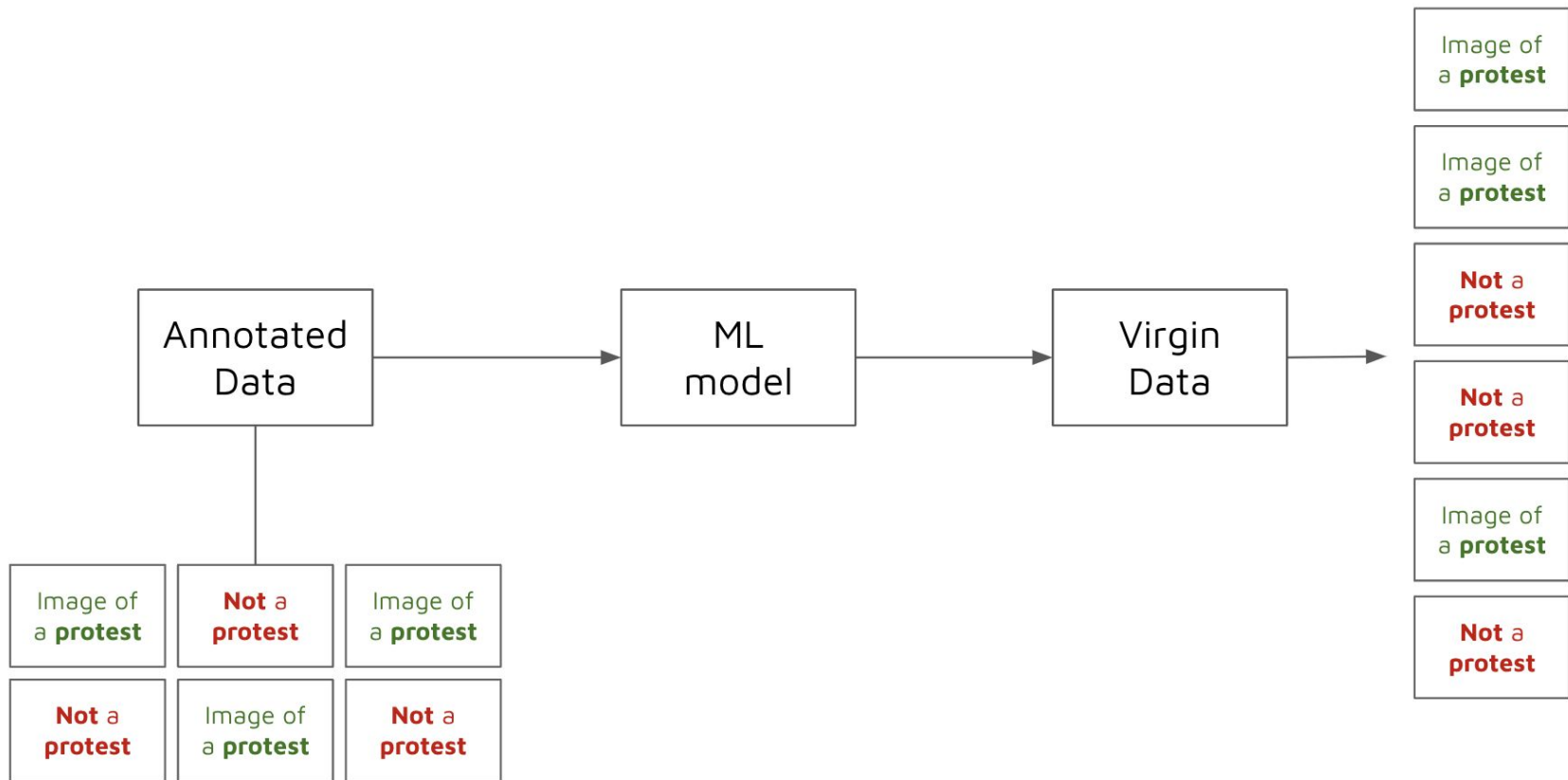
(Political) scientists have an **increasing amount of (digital) data available** for addressing their questions of interest



The volume of data generated, consumed, copied, and stored is projected to exceed 180 zettabytes by 2025



Supervised machine learning is often used to identify a defined concept in large amounts of data



Often this data is **multimodal**: with text, visuals, audio



Donald J. Trump  @realDonaldTrump · Sep 21

Hello everyone! I have something incredible to share today, as we are introducing the launch of our Official Trump Coins! The ONLY OFFICIAL coin designed by me—and proudly minted here in the U.S.A. The President Donald J. Trump First Edition Silver Medallion will be available

[Show more](#)



 24K

 27K

 144K

 26M

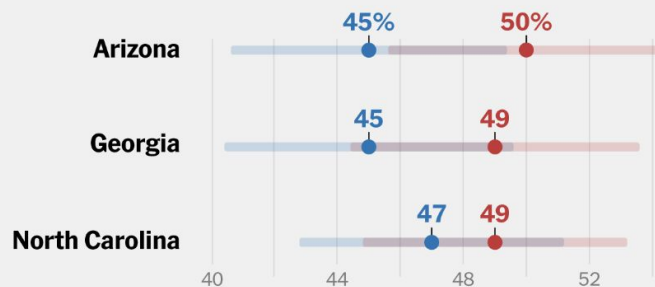
 

THE NEW YORK TIMES/SIENA COLLEGE POLL

Sept. 17 to 21


If the 2024 presidential election were held today,
who would you vote for if the candidates were

Kamala Harris and **Donald Trump**?



Among likely voters. Shaded areas represent margins of error.

Yet **we mostly only use one data modality** to train supervised ML (e.g. text)

 **Donald J. Trump**  @realDonaldTrump · Sep 21

Hello everyone! I have something incredible to share today, as we are introducing the launch of our Official Trump Coins! The ONLY OFFICIAL coin designed by me—and proudly minted here in the U.S.A. The President Donald J. Trump First Edition Silver Medallion will be available

[Show more](#)



24K

27K

144K

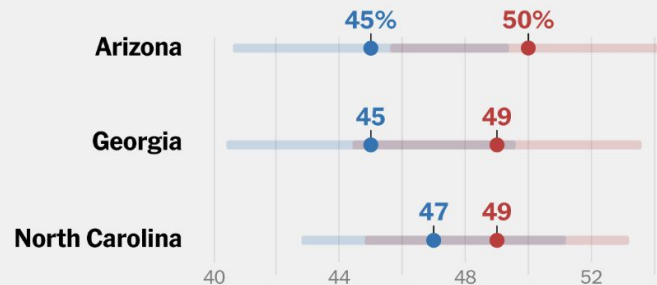
26M

THE NEW YORK TIMES/SIENA COLLEGE POLL

Sept. 17 to 21

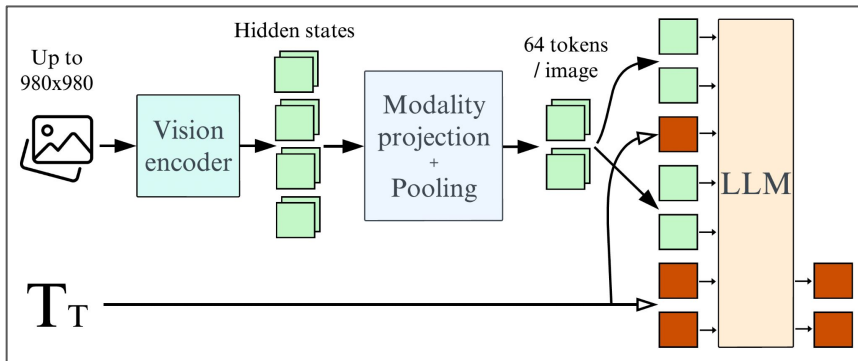
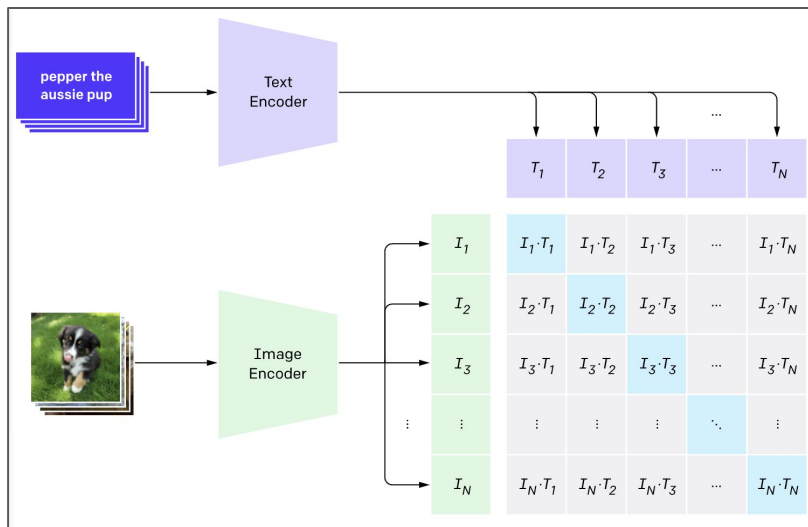
If the 2024 presidential election were held today,
who would you vote for if the candidates were
Kamala Harris and **Donald Trump**?



Among likely voters. Shaded areas represent margins of error.

Recent computational advances make **multimodal modeling** easier

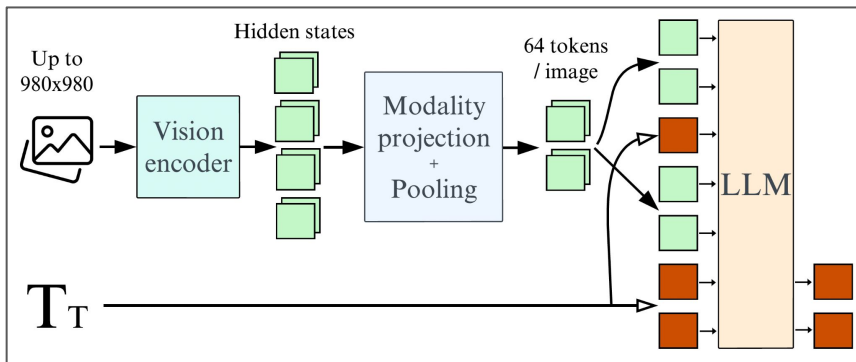
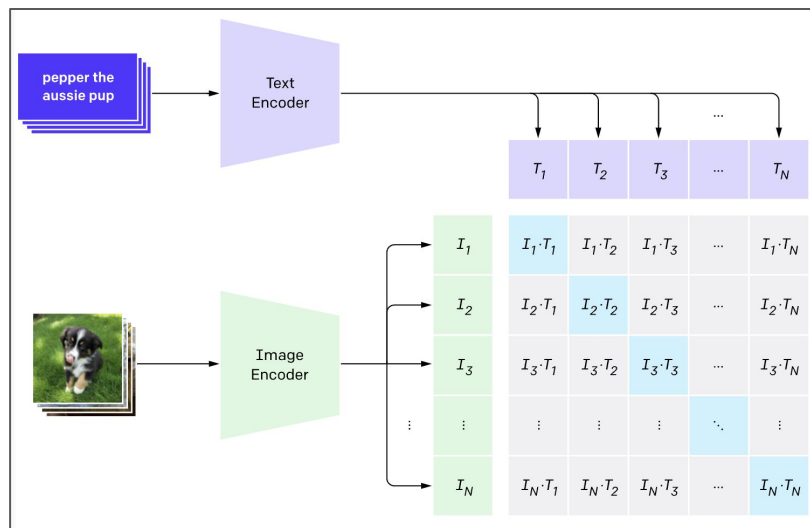
CLIP



Idex2

But we don't know much about **whether** nor **the conditions under which multimodality can help improve the performance of supervised ML**

CLIP



Idex2

Two **original annotated dataset** (10 tasks)

- (1) **YouTube Videos** from channels posting on US politics (N = ~4,000) → 6 tasks
- (2) **Twitter Posts** from interest groups from US, ES, DK, GE (N = ~4,000) → 4 tasks

Seven **models**:

- (1) **Text only**: SVM, BERT, Llama2, and Llama3
- (2) **Image only**: CNN
- (3) **Text + Image**: CLIP, Idefics2

This draft: 1 dataset and **4 tasks**

/!\ Work in Progress

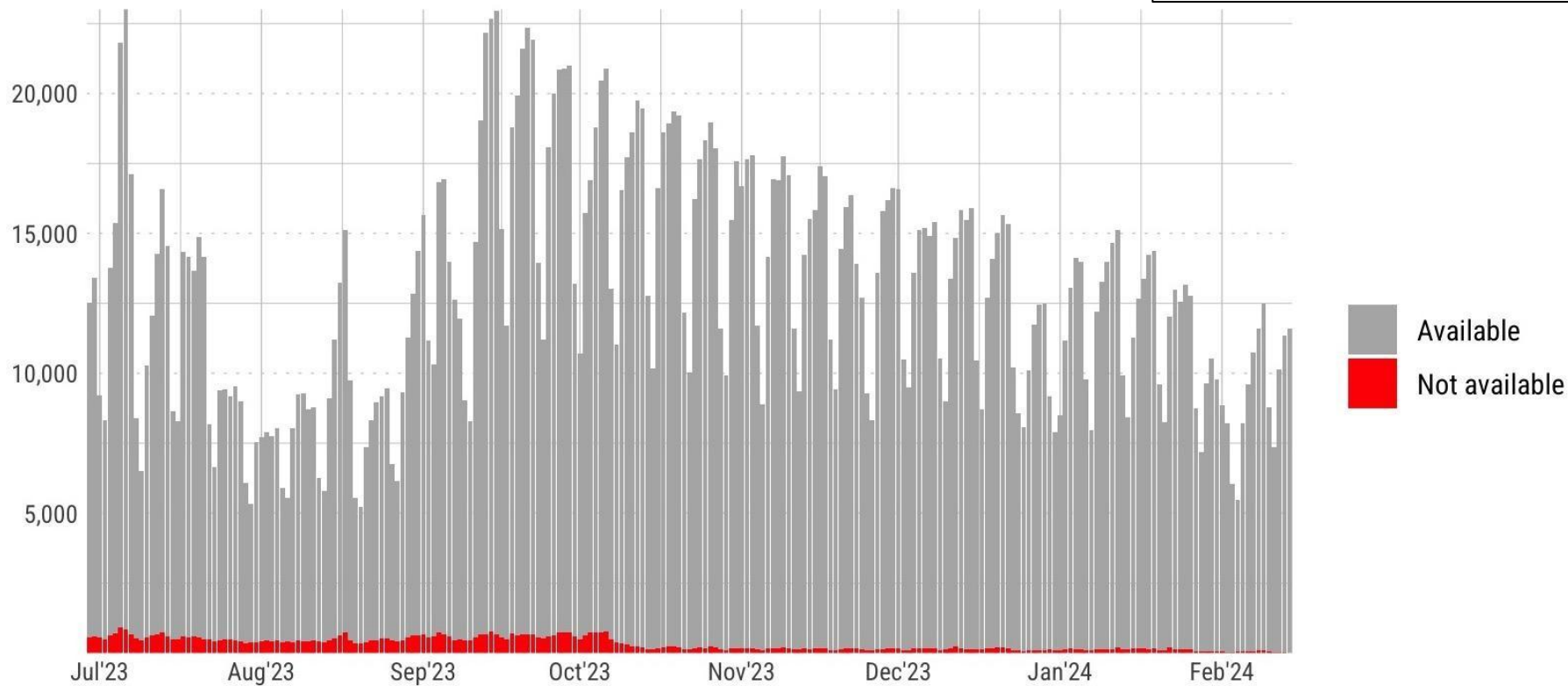
- (1) **YouTube Videos** from channels that post about US politics (N = ~4,000)
 - (a) *US politics*: is the video about US politics?
 - (b) *Hateful*: does the video have hateful content?
 - (c) *Typology*: what type of video? (e.g. opinion, high-quality news, etc.)
 - (d) *Ideology*: ideology of the video (liberal, moderate, conservative, neutral)

Seven **models**:

- (1) **Text only**: SVM, BERT, Llama2, and Llama3
- (2) **Image only**: CNN
- (3) **Text + Image**: CLIP, Idefics2

Data

- about **6 million videos**
- from about **12k channels**
- **3% channels** not available
- **2% videos** not available



Data

Metadata

- Title, Description, Topic labels, Duration, Tstamp

Video data

Transcript

- Tiny version
- Multilingual
- Quite fast (~1.5min. for 10min. video)



Ultimately, a Twitter bio deal he has said that then a vice President Biden in May 2017 and 2018 seemed to show that he met with then a vice President Biden in early 2017 along with it seems like Hunter Biden and Joe Biden. And... [The transcript continues with a detailed account of the investigation into Hunter Biden's activities, including his connections to the Biden family, his work at the Biden family's law firm, and his involvement in the Ukrainian energy sector. It mentions the discovery of a large cache of documents in the White House, the involvement of the FBI and the DOJ, and the eventual release of the documents to the public. The text is a verbatim transcription of a video, with some minor corrections for clarity.]

Frames

- Download video (temporarily)
- Extract frames (e.g. 1 frame/sec.)
- Frame embeddings (MoblieNetv3 small)
- Cosine similarity all pairs of frames
- identify "duplicate" frames (e.g. > 90%)
- Only keep "unique" frames



Data

Task	Description	Values	N	%
US Politics	Whether the video is about, or relevant to, US politics	0	1,935	49.8%
		1	1,945	50.2%
		<i>N</i>	3,880	100.0%
Hateful	Whether the video contains hateful language/behavior	0	3,431	88.4%
		1	449	11.6%
		<i>N</i>	3,880	100%
Idology	The ideological leaning of the video	Neutral	238	21.7%
		Conservative	476	41.9%
		Moderate	176	15.5%
		Liberal	247	20.9%
		<i>N</i>	1,137	100%
Typology	Type of video	Campaign	16	1.4%
		Educational	61	5.2%
		Satire	73	6.2%
		Low-Qual News	108	9.2%
		High-Qual News	332	28.4%
		Opinion	581	49.6%
		<i>N</i>	1,171	100%

Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3

Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2

Models

Text only:

(1) SVM

(2) BERT

(3) Llama2

(4) Llama3

Image only

(5) CNN

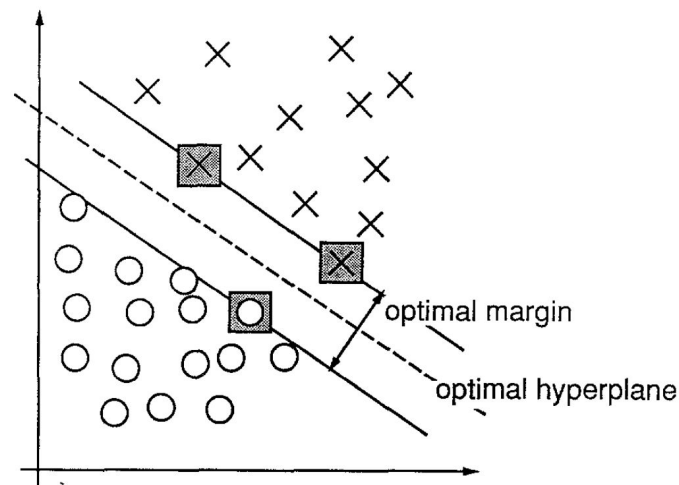
Text + Image

(6) CLIP

(7) Idefics2

- Ngram-based, **no prior language knowledge**, fully task dependent.

	Terms →					
	T1	T2	T3	...	TN	
D1	w11	w12	w13	...	w1N	
D2	w21	w22	w23	...	w2N	
D3	w31	w32	w33	...	w3N	
⋮	⋮	⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	⋮	⋮	
⋮	⋮	⋮	⋮	⋮	⋮	
DM	wM1	wM2	wM3	...	wMN	



Cortes & Vapnik (1995)

Models

Text only:

(1) SVM

(2) **BERT**

(3) Llama2

(4) Llama3

- Transformed-based, self-supervised, language model (prior knowledge), **fine-tuned on next sentence**, can be fine-tuned to do new tasks.
- Trained on: 11k books (800mil tokens) + English wikipedia (**2.5 bil** tokens)
- **bert-base-uncased**: 110 mil parameters
- Context length: **512** tokens

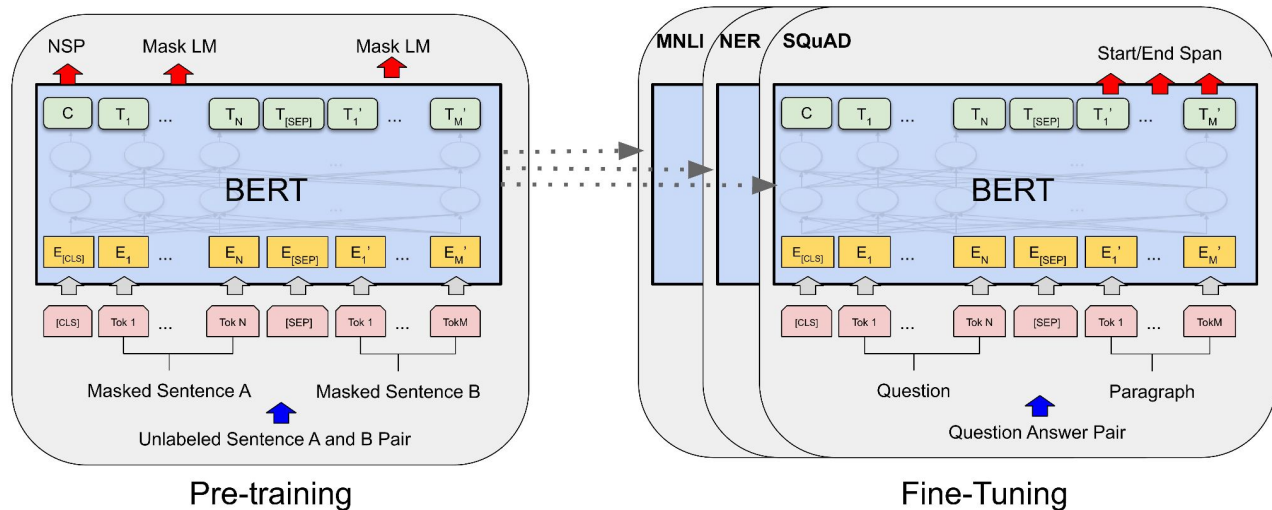
Image only

(5) CNN

Text + Image

(6) CLIP

(7) Idefics2



Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2**
- (4) Llama3

- Transformed-based, self-supervised, language model (prior knowledge), **fine-tuned on instruction task**, can be fine-tuned on new instructions
- Trained on: **2 tril** tokens, publicly available sources (although *unknown*)
- **7Bil/13bil/70bil** parameters → **Llama-2-7b-chat**
- Context length: **4,096 tokens**

Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2

Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3**

- Transformed-based, self-supervised, language model (prior knowledge), **fine-tuned on instruction task**, can be fine-tuned on new instructions
- Trained on: **15 tril** tokens, publicly available sources (although *unknown*)
- **8Bil/70/405 bil** parameters → **Llama-3-8b-instruct**
- Context length: **8,000**

Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2

Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3

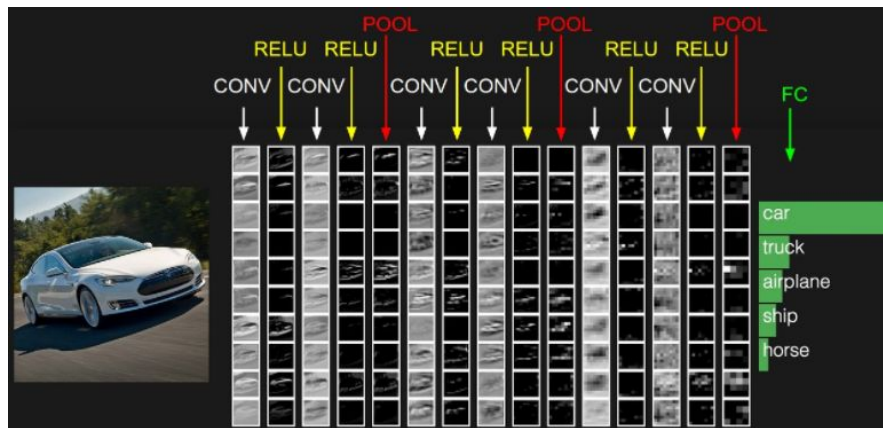
- Pre-trained for object recognition: 1,000 ImageNet object classes
- Trained on: **1.28 mil** images
- **25.6 mil** parameters → **ResNet50**
- Input size: 224 x 224 x 3

Image only

- (5) **CNN**

Text + Image

- (6) CLIP
- (7) Idefics2



Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3

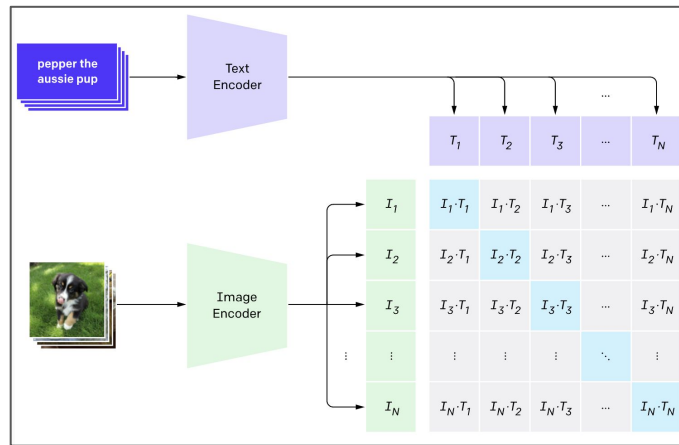
Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2

- Text and image encoder
- Trained on **400 mil** text-image pairs: e.g. image and its caption
- **Self-trained**: predicting correct text-image pair
- Image input size: **224 x 224 x 3**
- **150/400 mil** parameters → **ViT-B/32**
- Context length: **77 tokens**



Models

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3

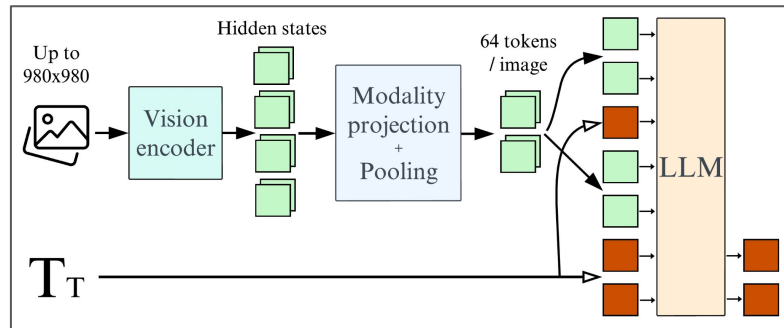
Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2**

- Text and image **transformer** encoder
- Trained on:
 - **interleaved image-text document: 350 mil images and 115 bil text tokens.**
 - **Image-text pairs**
 - **PDF OCR extraction: 40 mil**
 - **instruction/chat: 50 open-source datasets**
 - Total: **1.5 bil** images and **225 bil** text tokens
- **8 bil** parameters → **idefics2-8b**
- Image input size: **native resolution (up to 980 x 980)**
- Context length: ?



Set Up

Text only:

- (1) SVM
- (2) BERT
- (3) Llama2
- (4) Llama3

Image only

- (5) CNN

Text + Image

- (6) CLIP
- (7) Idefics2

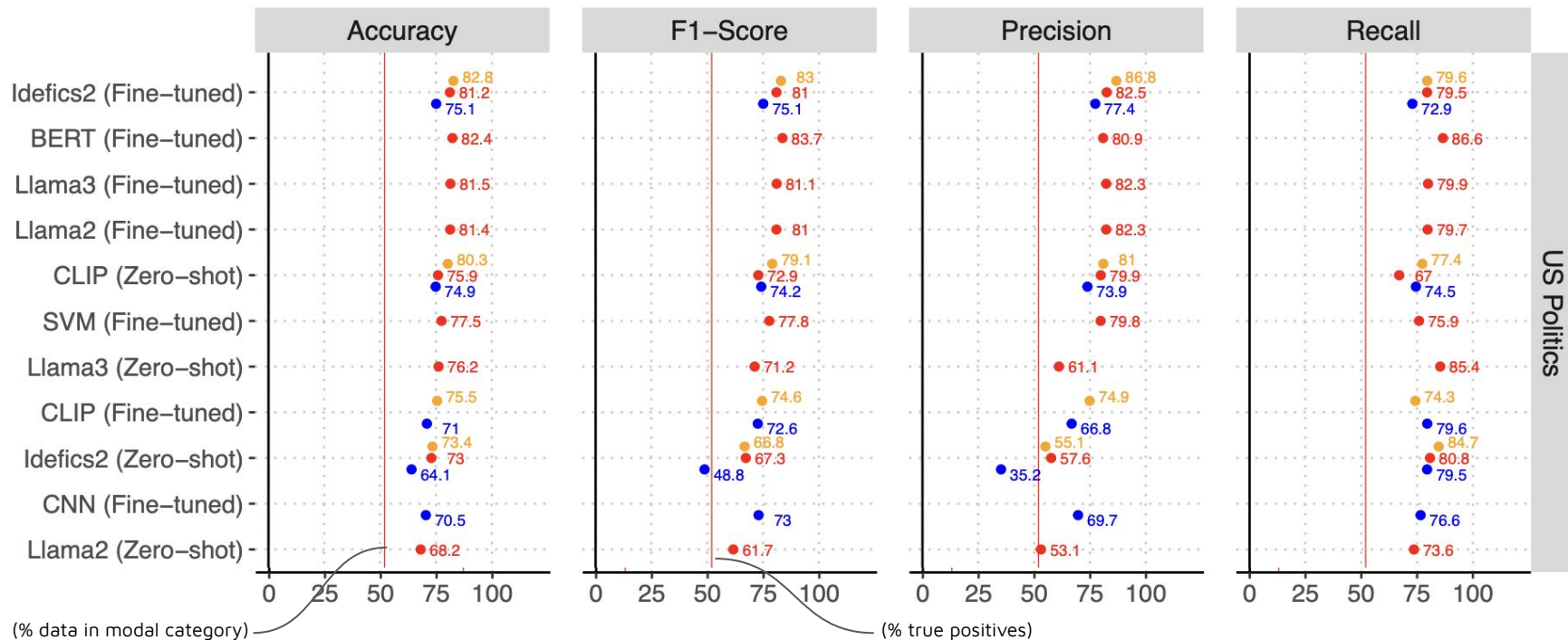
- Same **train** (80%) - **test** (20%) sets across models
- Further split train set 80/20 into train/**validation** per **fold**
- **3 folds** (today results only for 1 fold) and **10 epochs**/fold
- Image processing/input:
 - Resize + center_crop all images: 224 x 224 x3
- Text processing/input:
 - 20,000 token vocabulary; no stopwords; linear kernel
 - 512 tokens/video
 - 800 tokens/video
 - 77 tokens/video
- Fine-tuning:
 - $\text{Transcript}_v + \text{Frame}_{v,f} \rightarrow \text{Label}_v$
 - BERT, CNN: new prediction head for each of our tasks
 - Llama2, Llama3, Idefics: same prompts across models
- Evaluation of image and text + image models:
 - Binary: $\text{Pred} = 1$ if $\text{sum}(\text{Frame}_{v,f}) > \text{threshold}$ (=0 otherwise)
 - Multiclass: $\text{Pred} = \text{mode}(\text{Frame}_{v,f})$

A quick look at how to fine tune a VLM (Idefics2)

```
messages = [  
  {  
    "role": "user",  
    "content": [  
      {"type": "image"}, # frame_v1_f1  
      {"type": "image"}, # frame_v1_f2  
      {"type": "image"}, # frame_v1_f3  
      {"type": "text", "text": transcript_v1,  
      {"type": "text", "text": "Is the previous text and images about or relevant to US politics? Answer YES or NO."}  
    ],  
  },  
  {  
    "role": "assistant",  
    "content": [  
      {"type": "text", "text": "YES"}  
    ]  
  }  
]
```

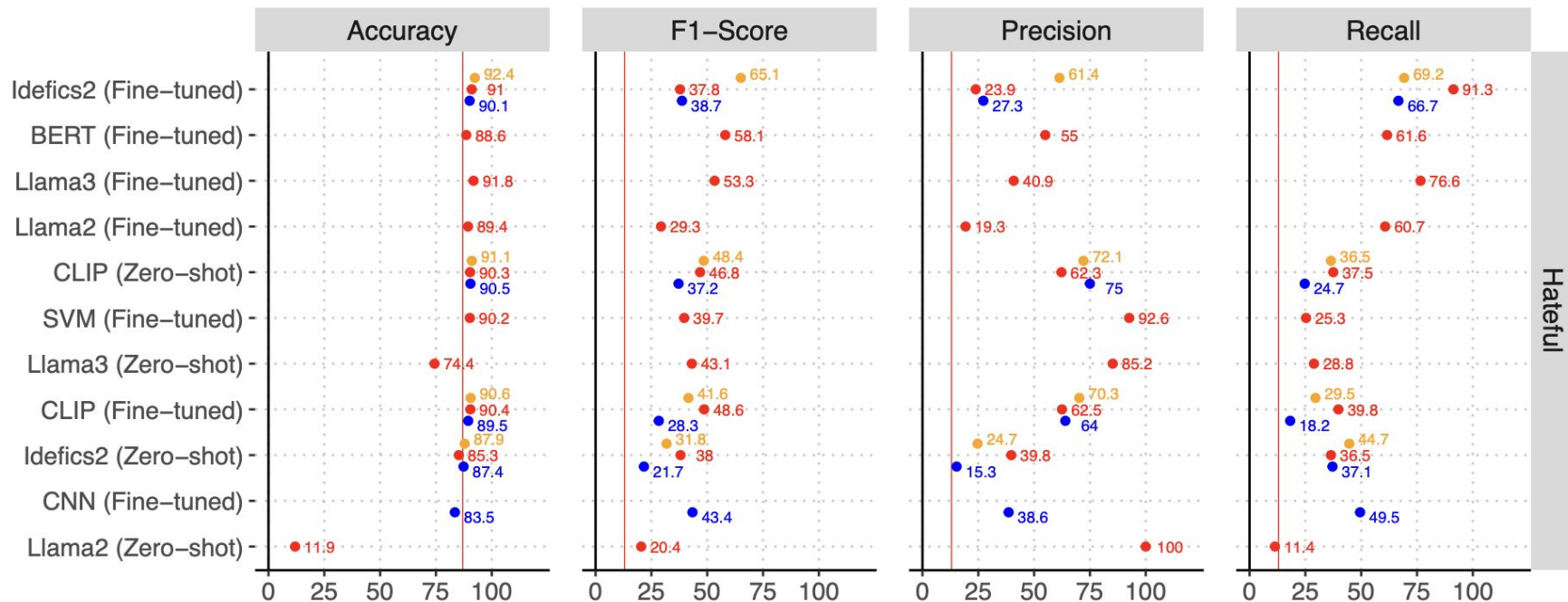
Results: overview

● Image ● Multimodal ● Text



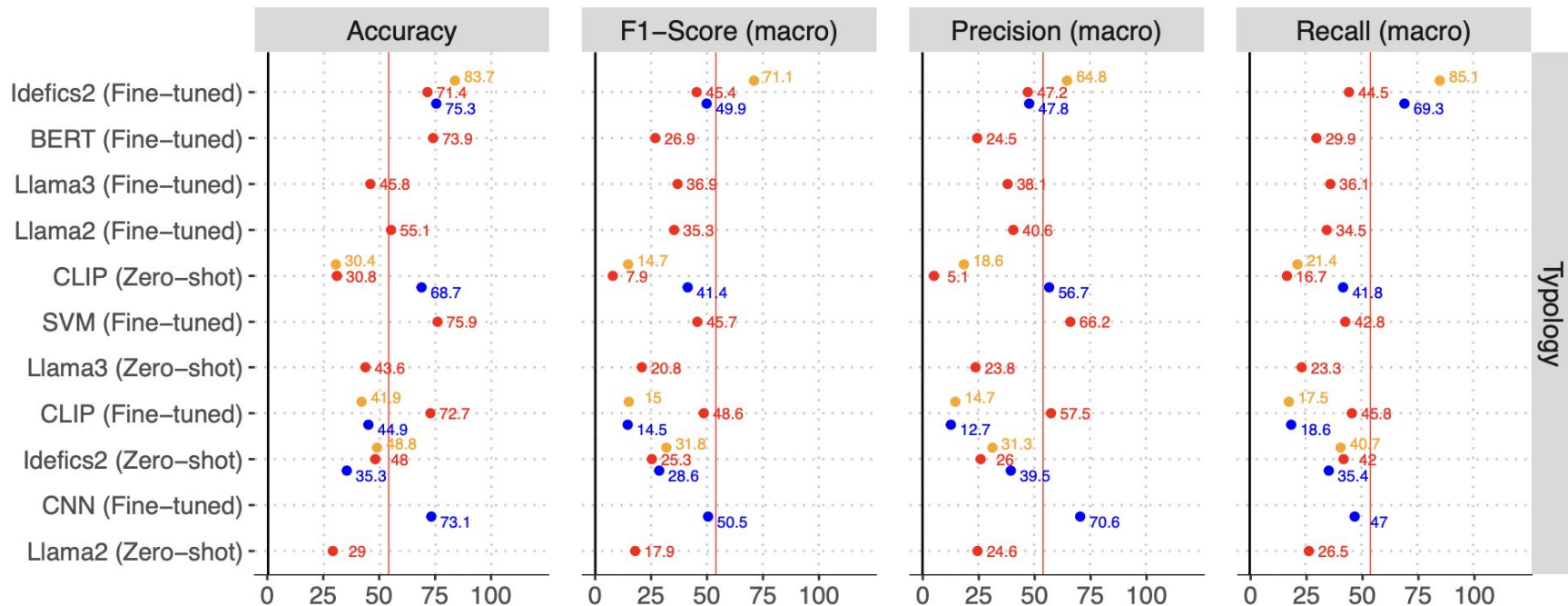
Results: overview

• Image • Multimodal • Text



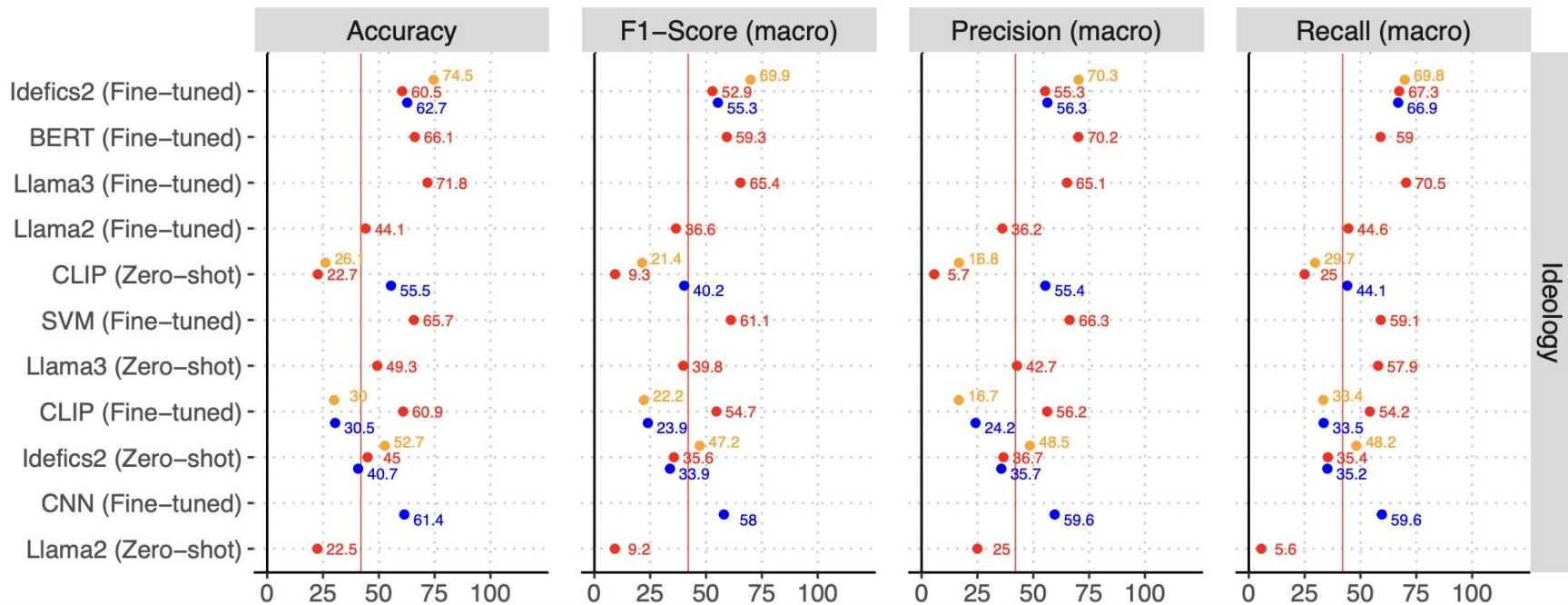
Results: overview

• Image • Multimodal • Text

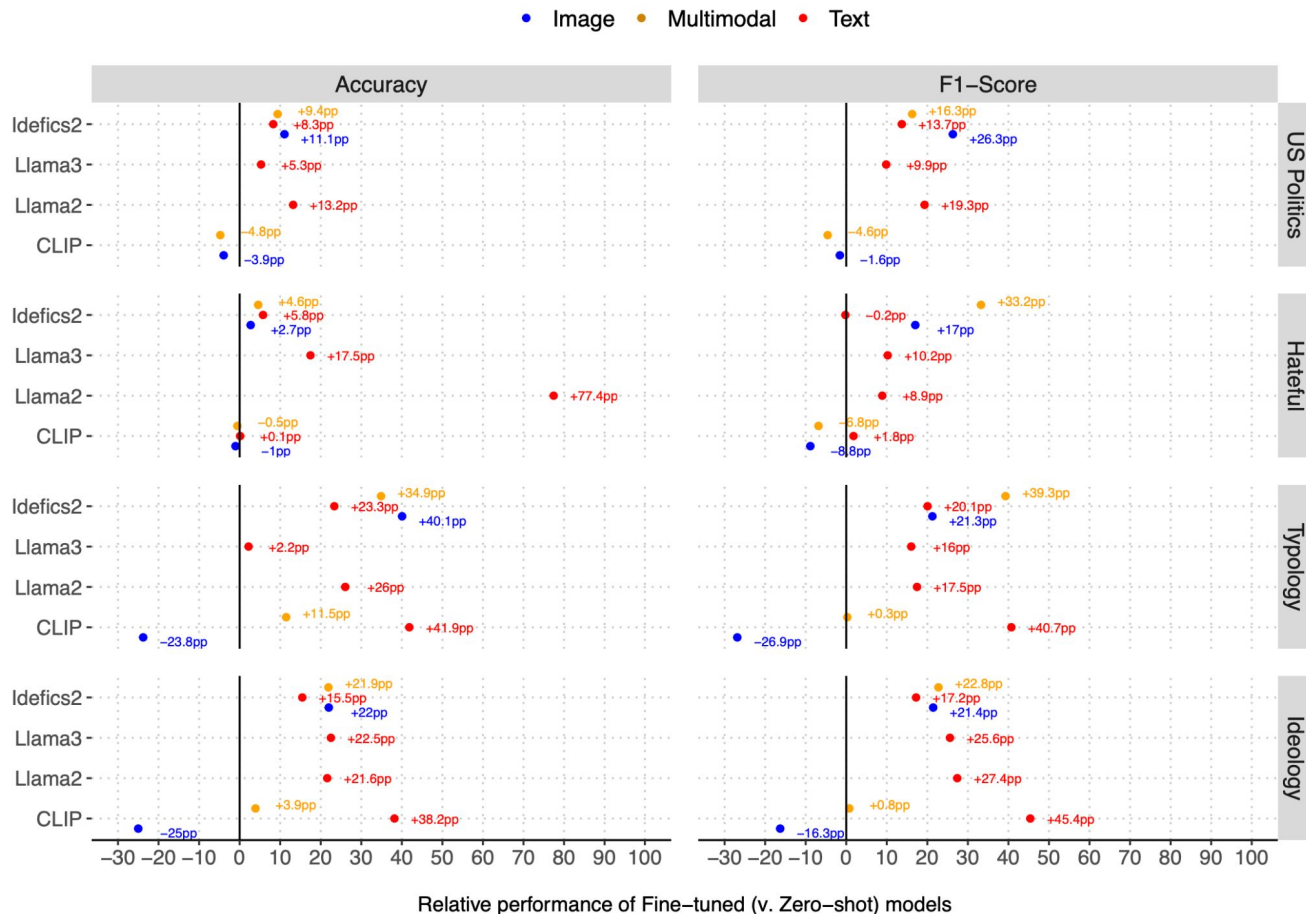


Results: overview

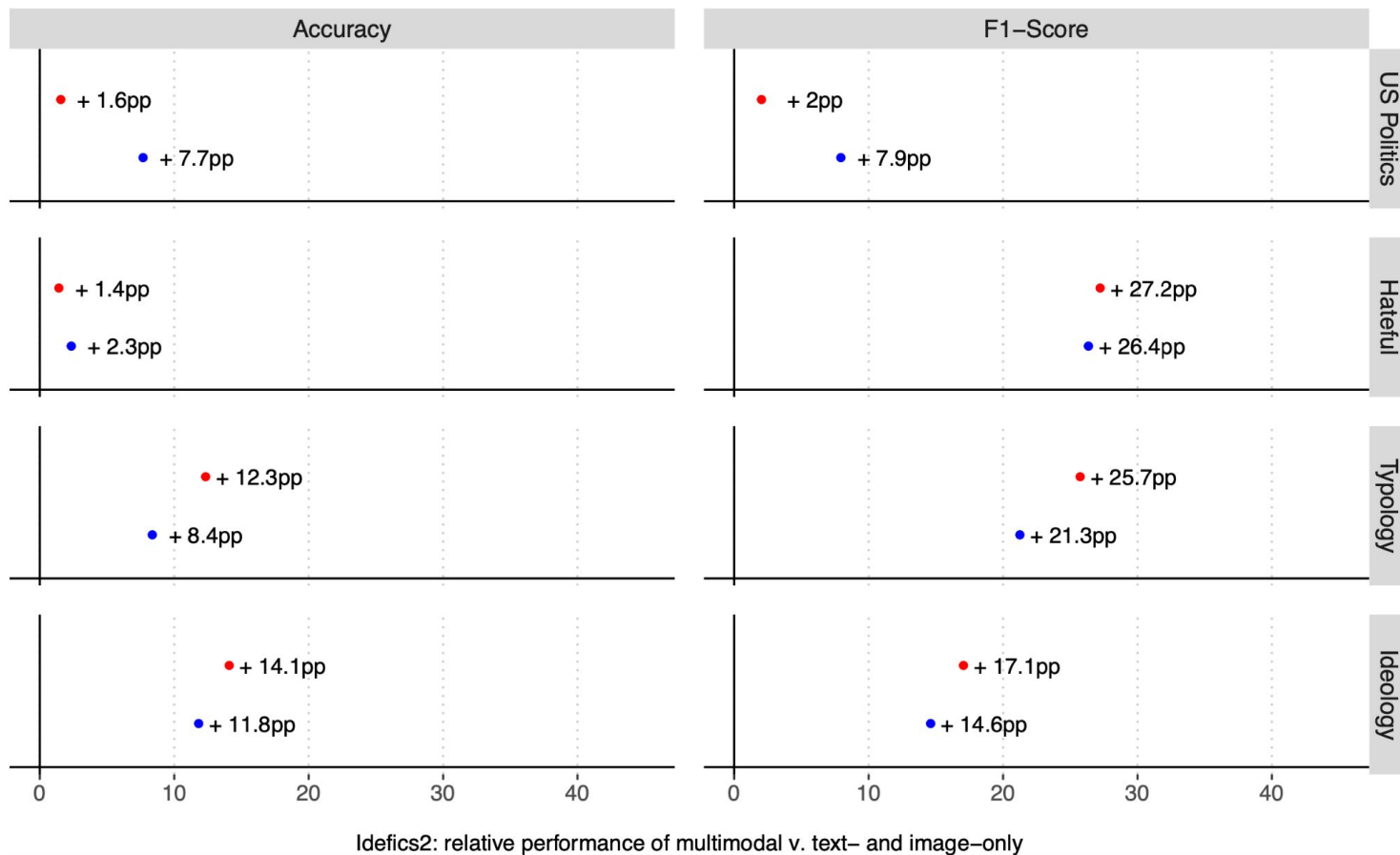
• Image • Multimodal • Text



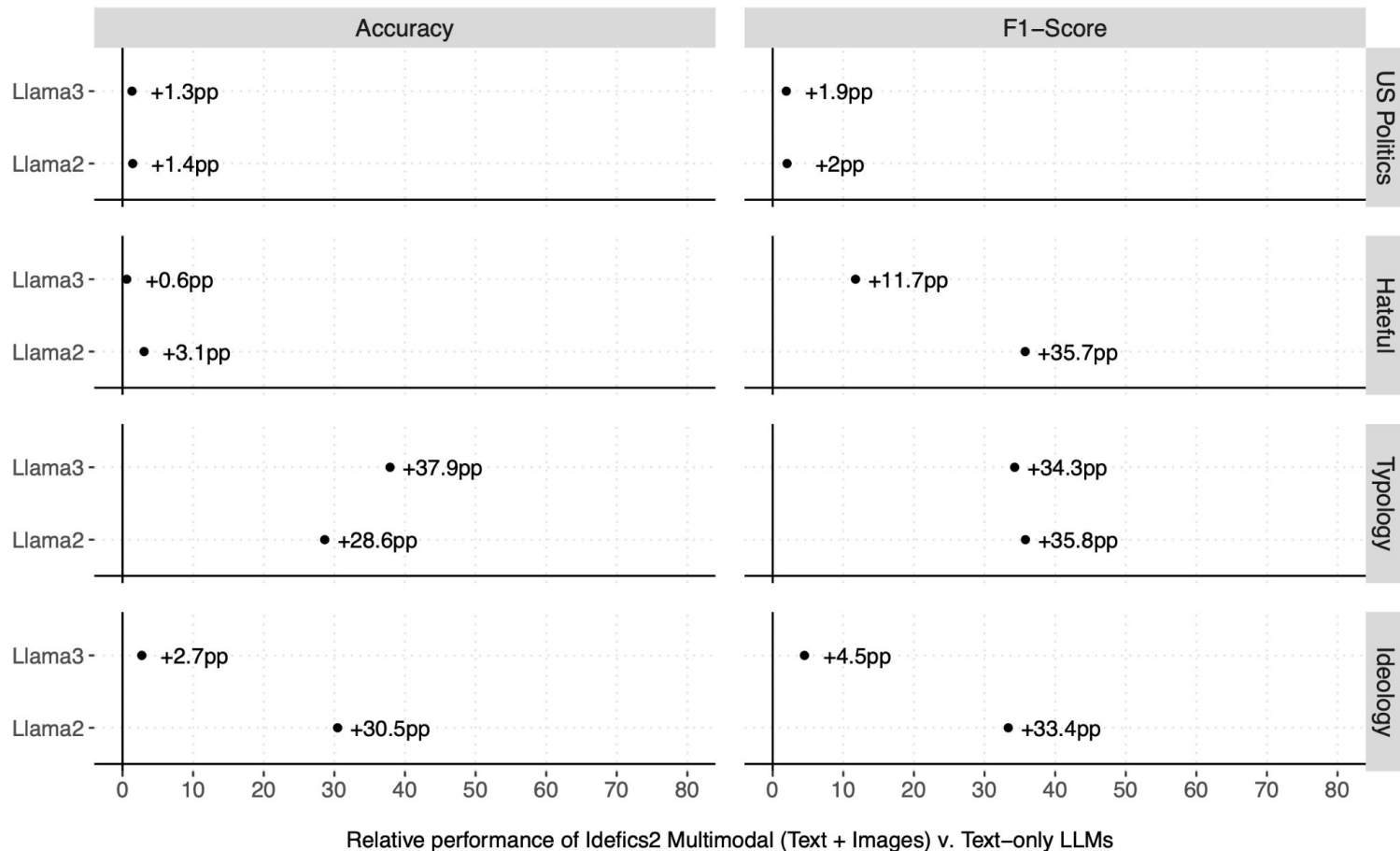
Results: fine-tuning makes a big difference



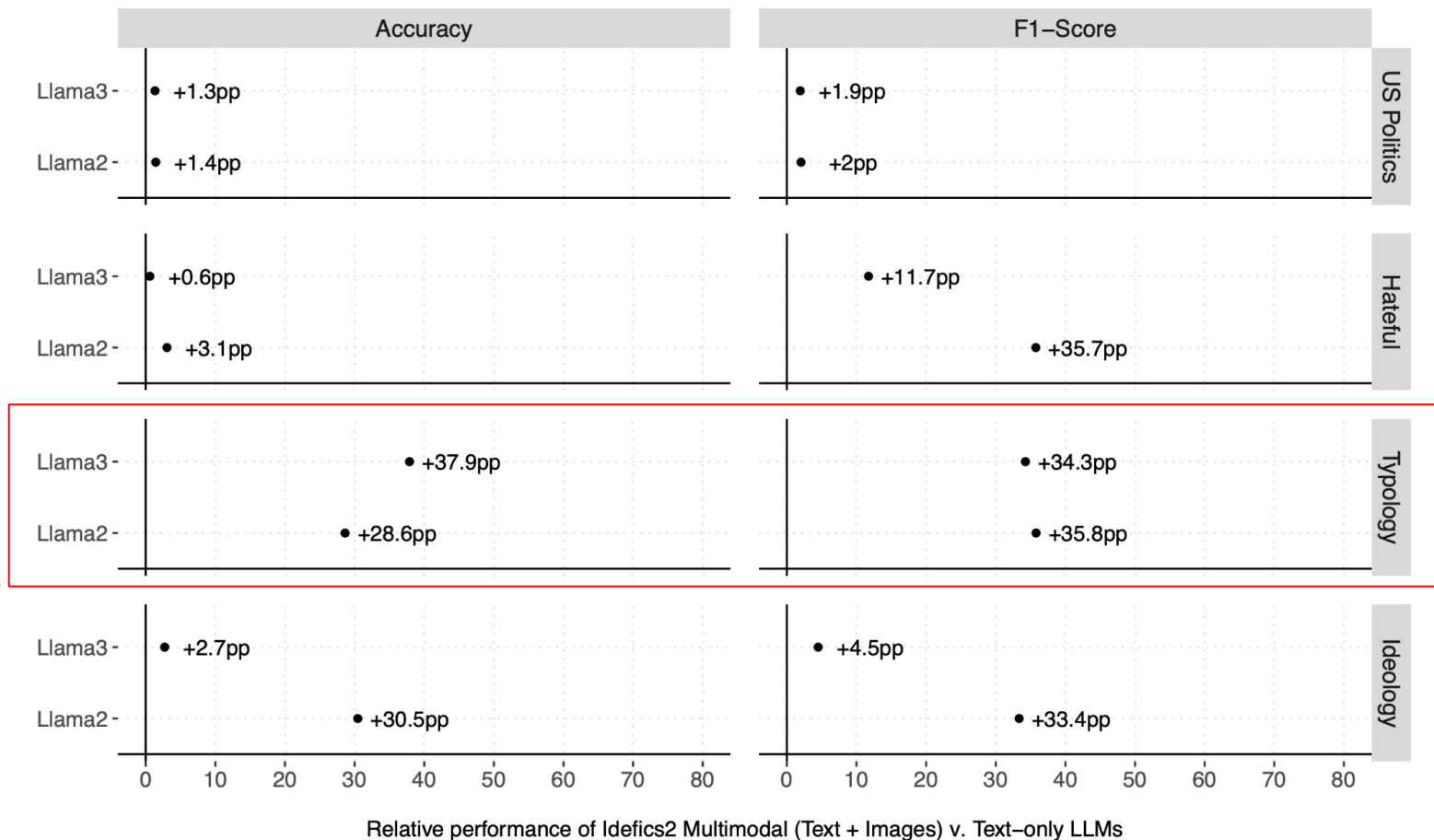
Results: VLM performs better in multimodal settings



Results: VLM outperforms SOTA open-source LLMs



Results: ... particularly on more visual-dependent task



Limitations

- Current results based on only 1 fold
- Smaller Llama2 and Llama3 models (7B v. 40-400B parameters)
- Only max. of 800 text tokens per video
- Probably some more prompt engineering is needed
- NEW promising open-source VLM: Idefics3-8B-llama3

Conclusions and next steps

- Multimodality helps to improve performance (v. text-only LLMs)
- Particularly on more visual-dependent tasks
- Next steps:
 - address limitations discussed in previous slide (folds, prompts, add larger/new models)
 - run the same computational experiments on the missing target variables for the YouTube dataset; and also on the Interest Group Twitter data
 - potentially adding one more dataset (TikTok) with higher visual dependence