

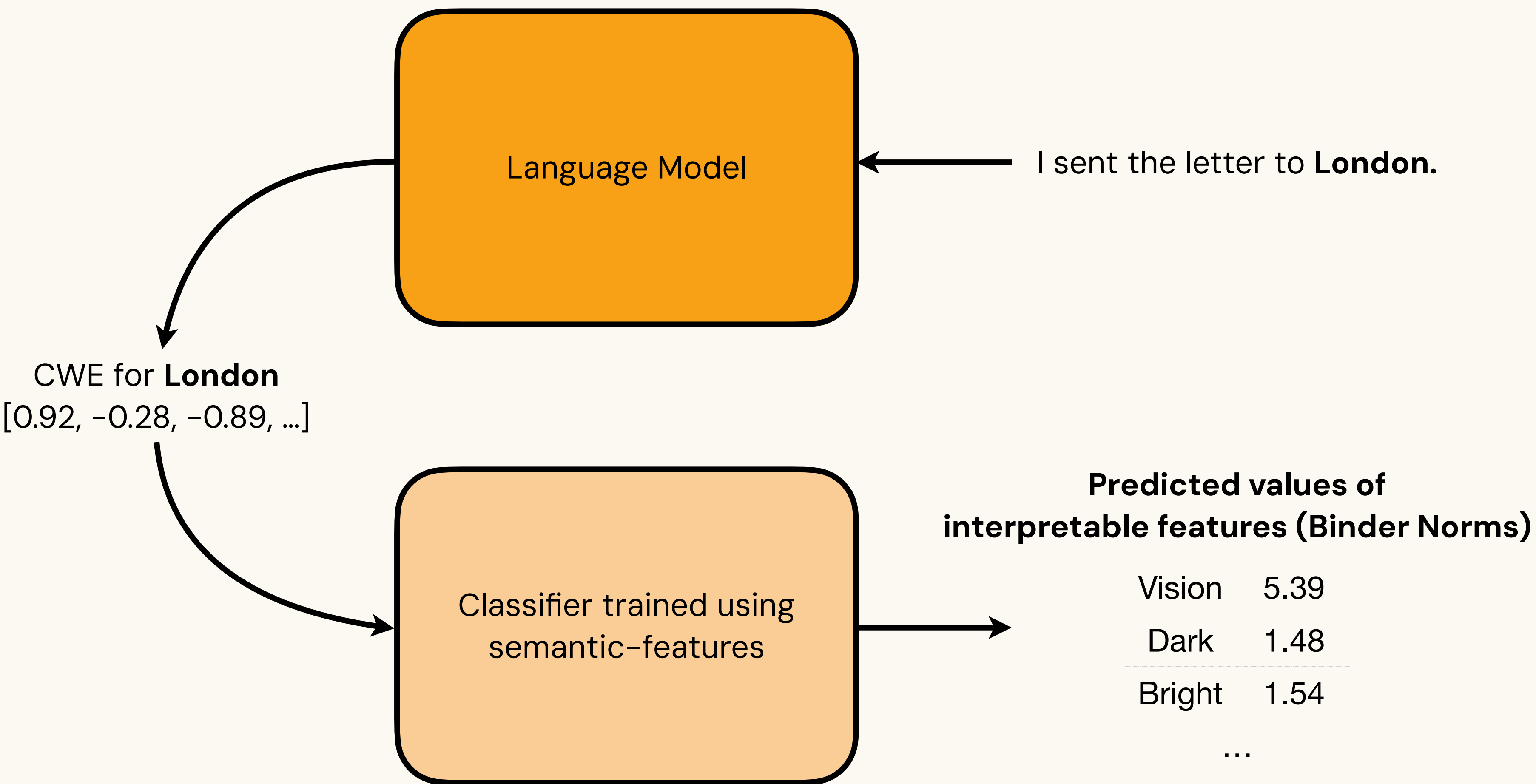
semantic-features: A User-Friendly Tool for Studying Contextual Word Embeddings in Interpretable Semantic Spaces



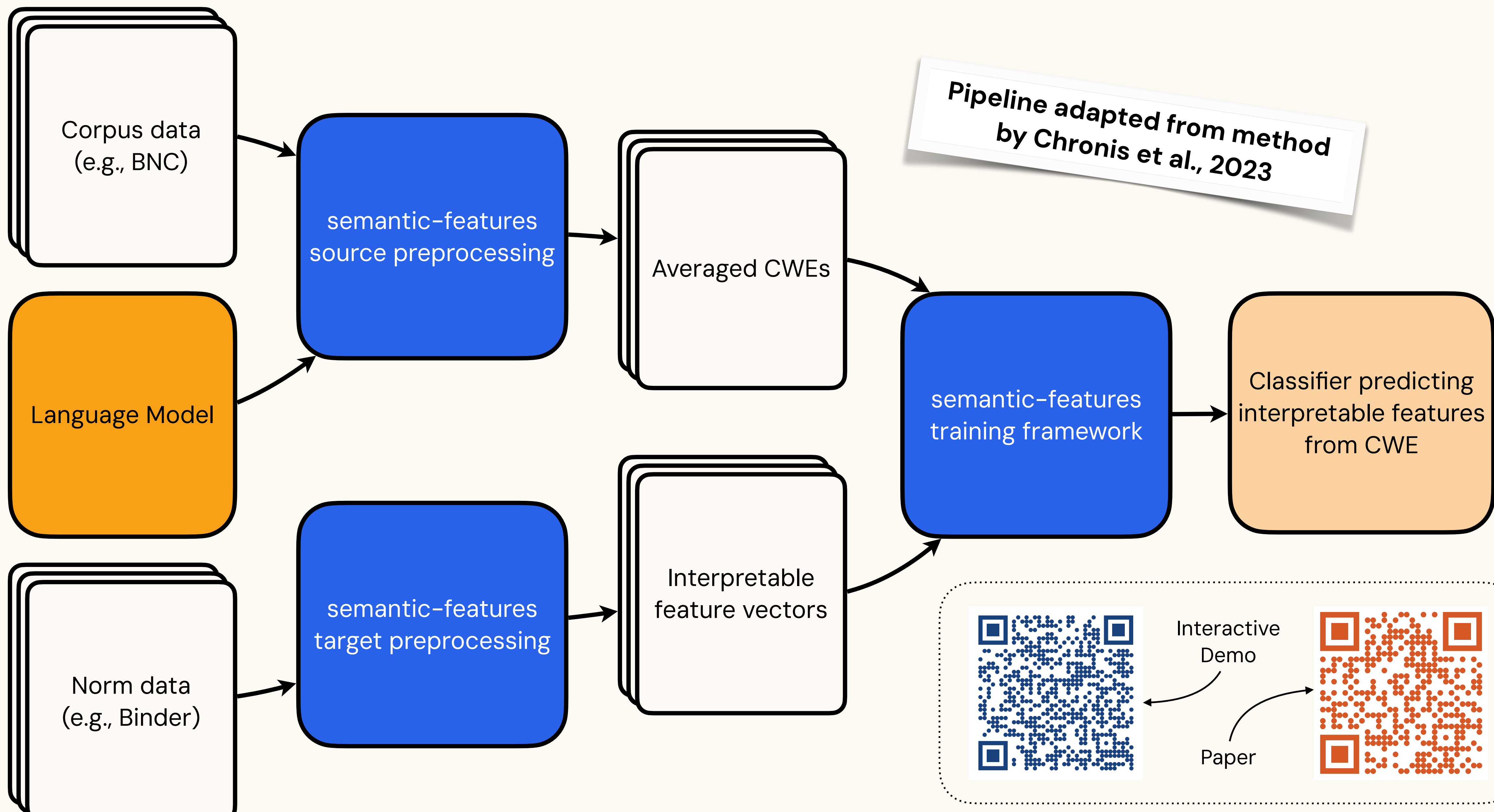
Jwalanthi Ranganathan¹, Rohan Jha¹, Kanishka Misra^{1*,2}, Kyle Mahowald¹

^{*} Work partly done as a Postdoc at The University of Texas at Austin

Prediction Backend



Classifier Training Using semantic-features



Predictions Using semantic-demo

Spaces

jwalanthi/semantic-demo

like 0

Running

...

Sentence

I sent the letter to London.

Word

London

LLM

☒ BERT ☐ ALBERT ☐ RoBERTa

Norm

☒ Binder ☐ McRae ☐ Buchanan

Layer

8

Clear

Submit

output

All Positive Predicted Values:

Vision	5.39
Scene	4.43
Large	3.96
Shape	3.95
Weight	3.92
Practice	3.6
Landmark	3.43
Benefit	3.37
Needs	2.72
Complexity	2.5
Attention	2.34
Near	2.33
Texture	2.25
Color	2.23
Pattern	2.2
Pleasant	2.07
Arousal	1.91
Touch	1.8
Social	1.66
Consequential	1.62
Path	1.61

Case Study: Recipient Semantics in Dative Constructions

Prepositional Object (PO):
I sent the letter to **London**.



Double Object (DO):
I sent **London** the letter.



(Goldberg, 1995; Hovav and Levin, 2008; Beavers, 2011)

Feature	Definition
Biomotion	showing movement like that of a living thing
Body	having human or human-like body parts
Human	having human or human-like intentions, plans, or goals
Face	having a human or human-like face
Speech	someone or something that talks
Landmark	having a fixed location, as on a map
Scene	bringing to mind a particular setting or physical location

Table 1: Feature definitions from Binder et al. (2016).

Feature	DO	PO
Biomotion	1.19	0.43
Body	1.00	0.26
Human	0.89	0.48
Face	0.71	0.19
Speech	0.68	0.13
Landmark	1.83	3.43
Scene	2.59	4.43

Table 2: Relevant Binder features predicted for “London” in (1) using CWEs from BERT layer 8. The PO construction lends itself more towards “location” features, and the DO more towards animate features.



Figure 2: For each layer of an LM, we extract the CWE and project it into Binder space. **Left:** we measure the average change across the test sentences in Person features from PO to DO. The positive values indicate that recipients in the DO are found to be more animate. **Right:** we measure the average change across test sentences in Place features from DO to PO. Here, the positive values indicate that the recipients were found to be more place-like in the PO.