

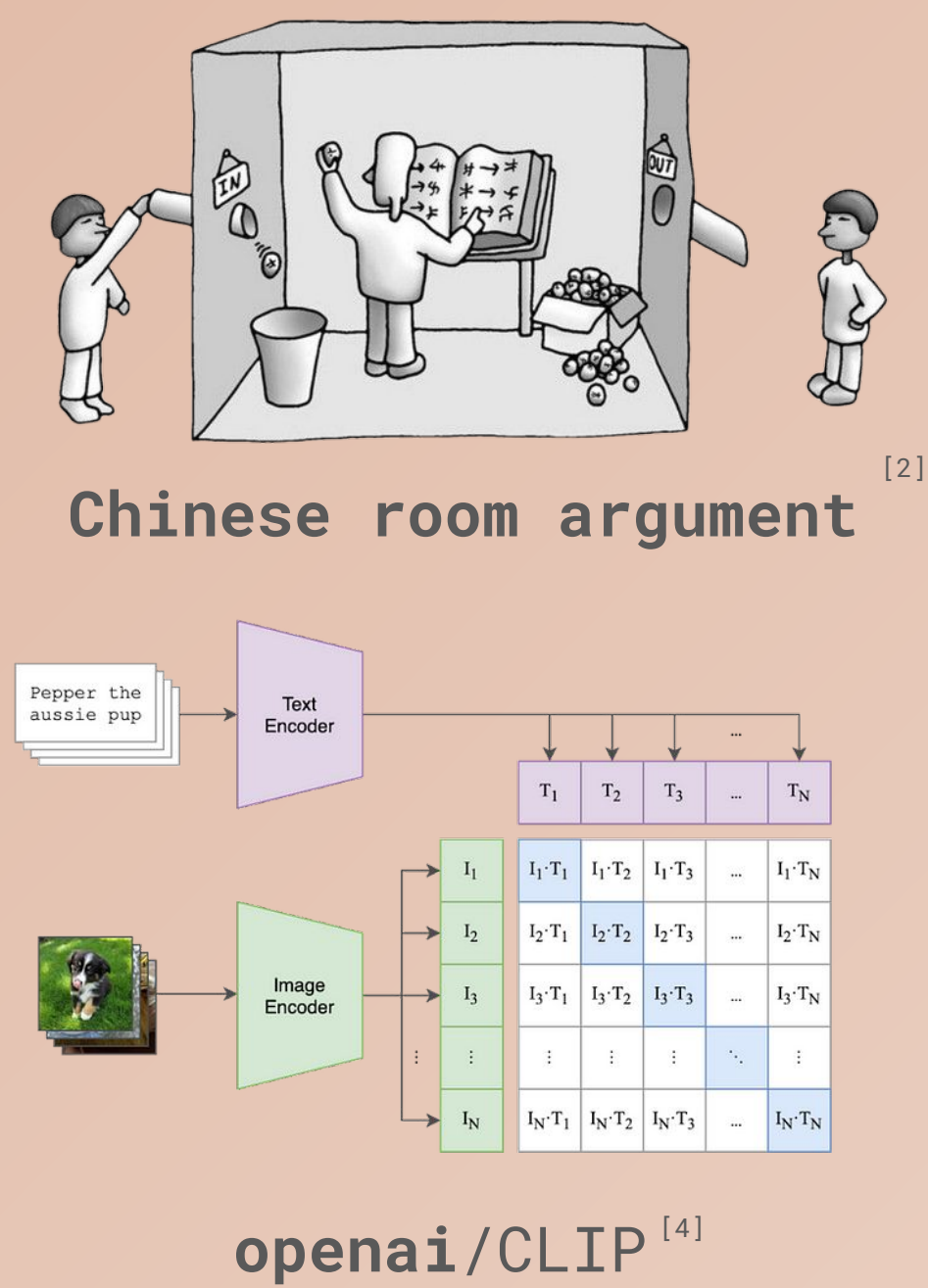
Exploring Visually Grounded Word Embeddings

Milan Miletić, Ryan Ott, Adrian Sauter

Context

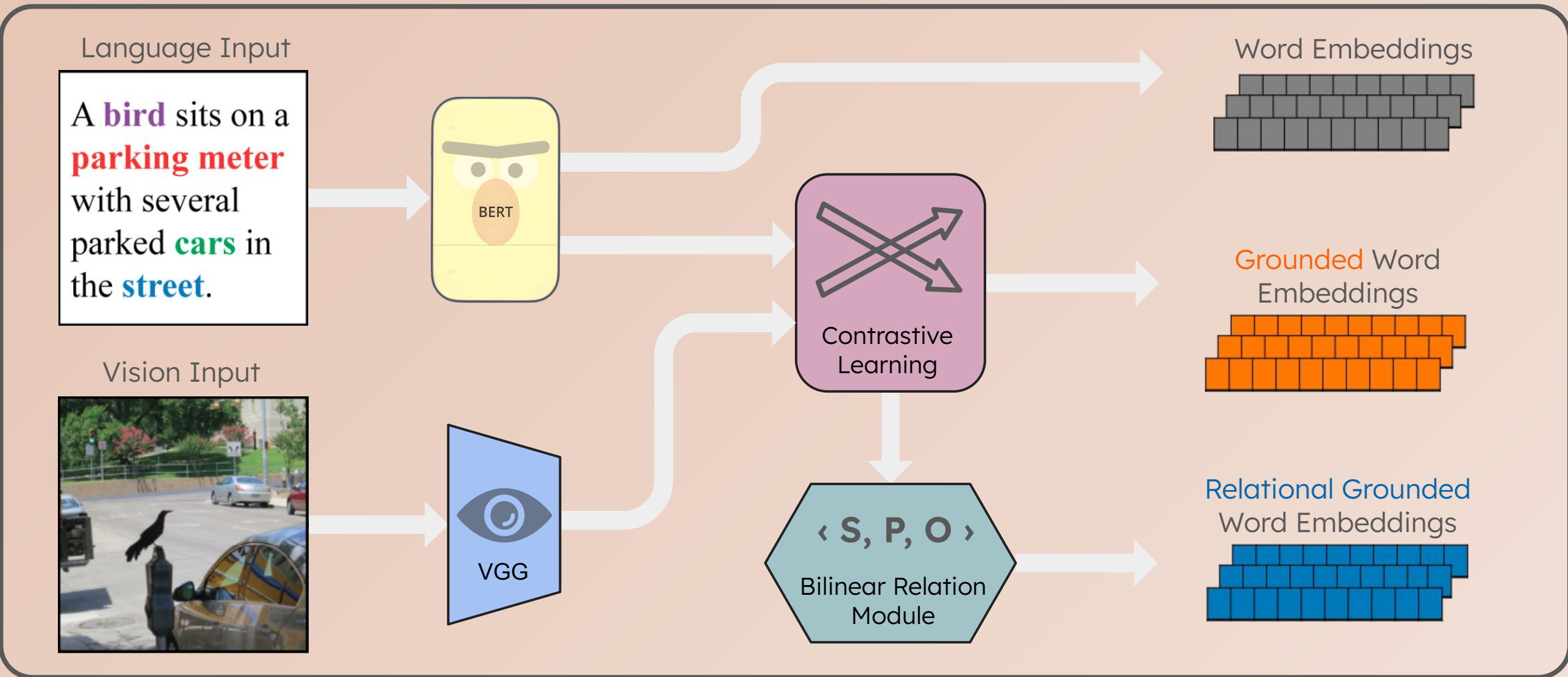
RED
YELLOW
GREEN
BLUE

Stroop task [3]



"Explainable Semantic Space by Grounding Language to Vision with Cross-Modal Contrastive Learning" (Zhang et al., 2021) [1]

Paper



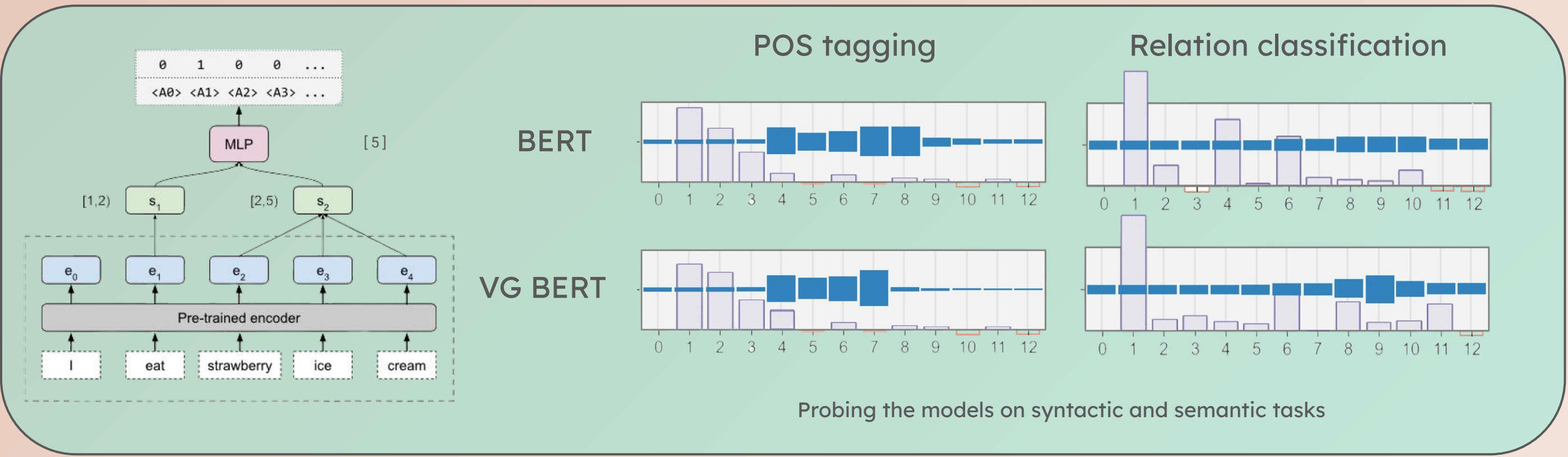
Research Questions

1. How is information propagated through the VG encoder?
2. Does VG equally impact concrete and abstract concepts?
3. Does VG help with resolving lexical ambiguities?
4. Are the grounded embeddings clustered better?

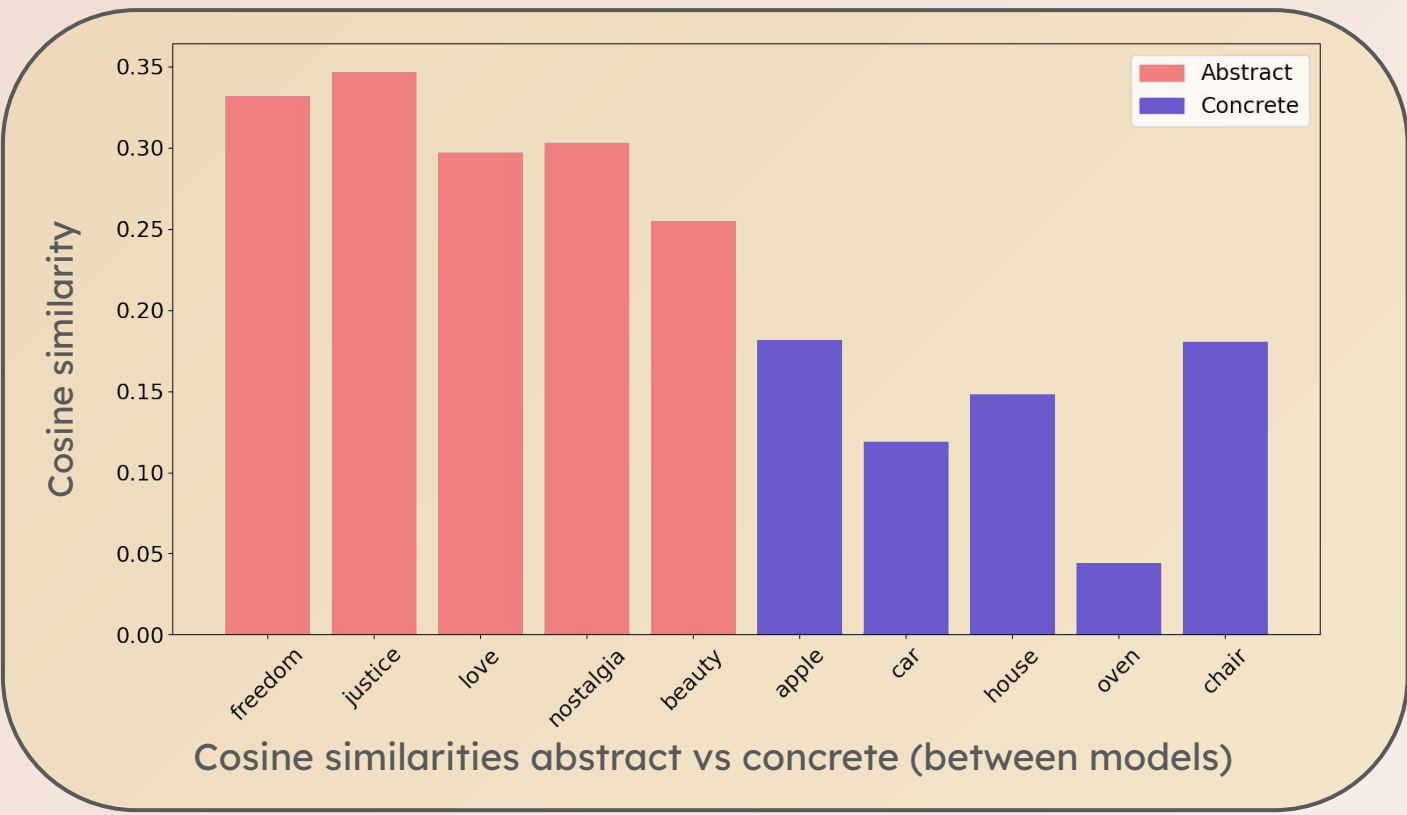
Results

	Probing		Lexical Ambiguity	Clustering
	POS	Rel. \uparrow	Cosine sim. \downarrow	Silhouette coef. \uparrow
BERT	0.96	0.49	0.56 ± 0.13	0.02
VG BERT	0.96	0.52	0.34 ± 0.15	0.31

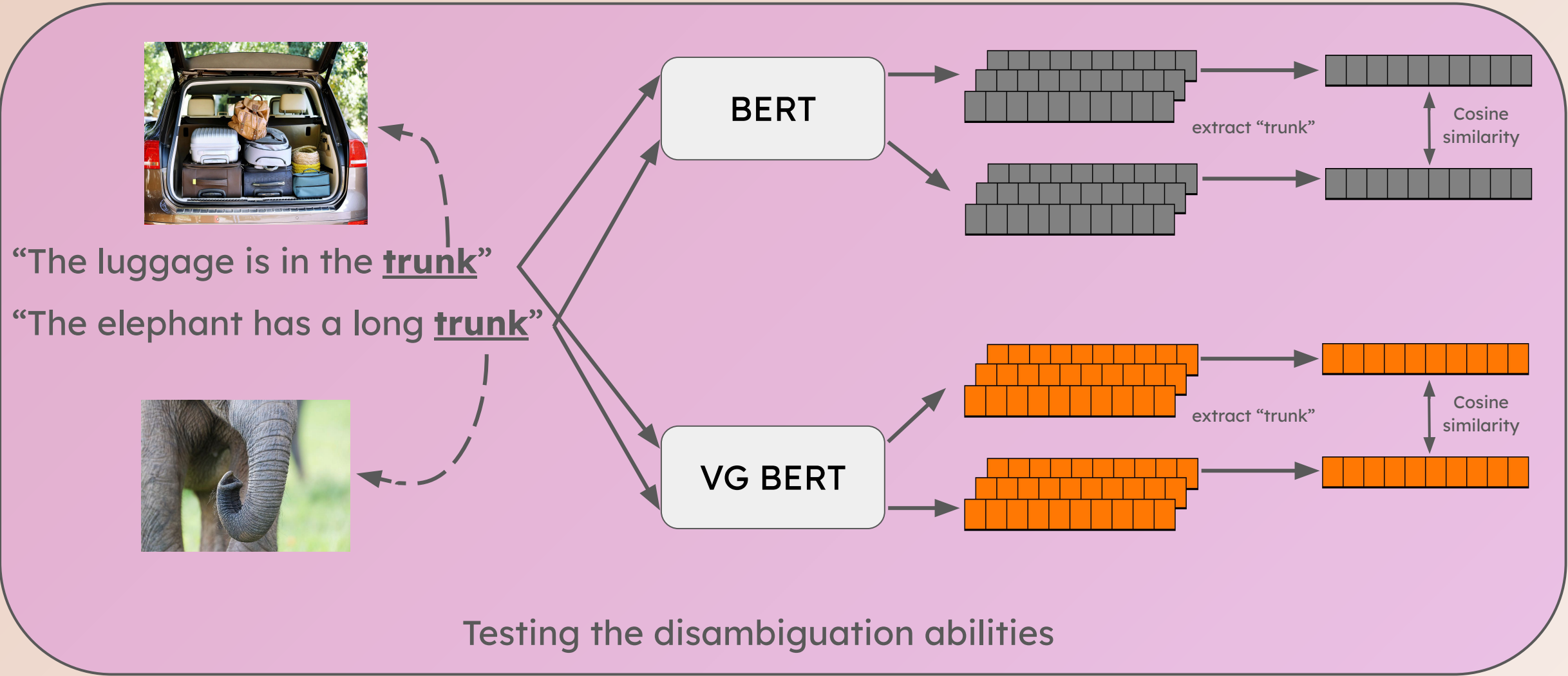
Probing



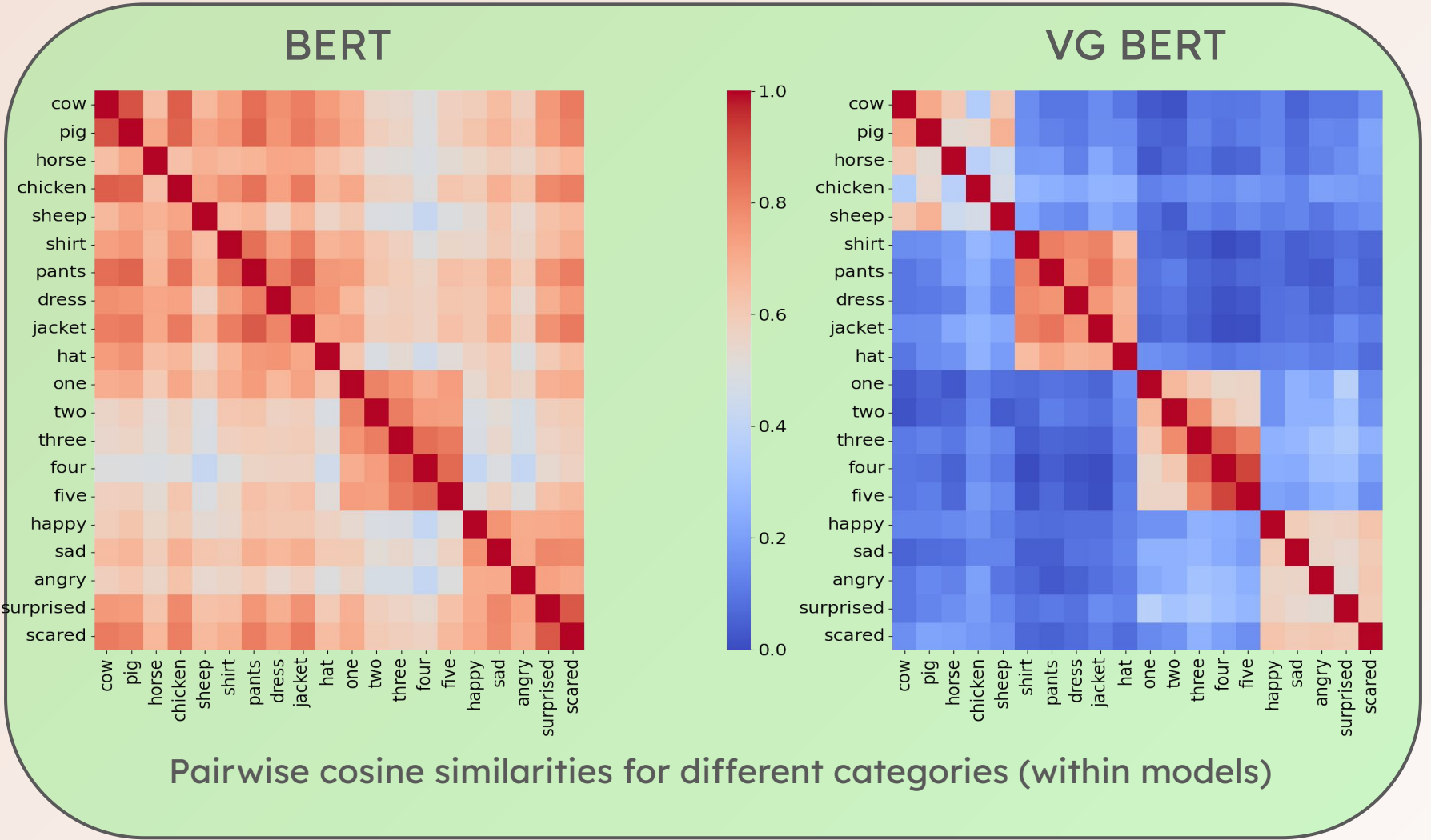
Abstract vs. Concrete



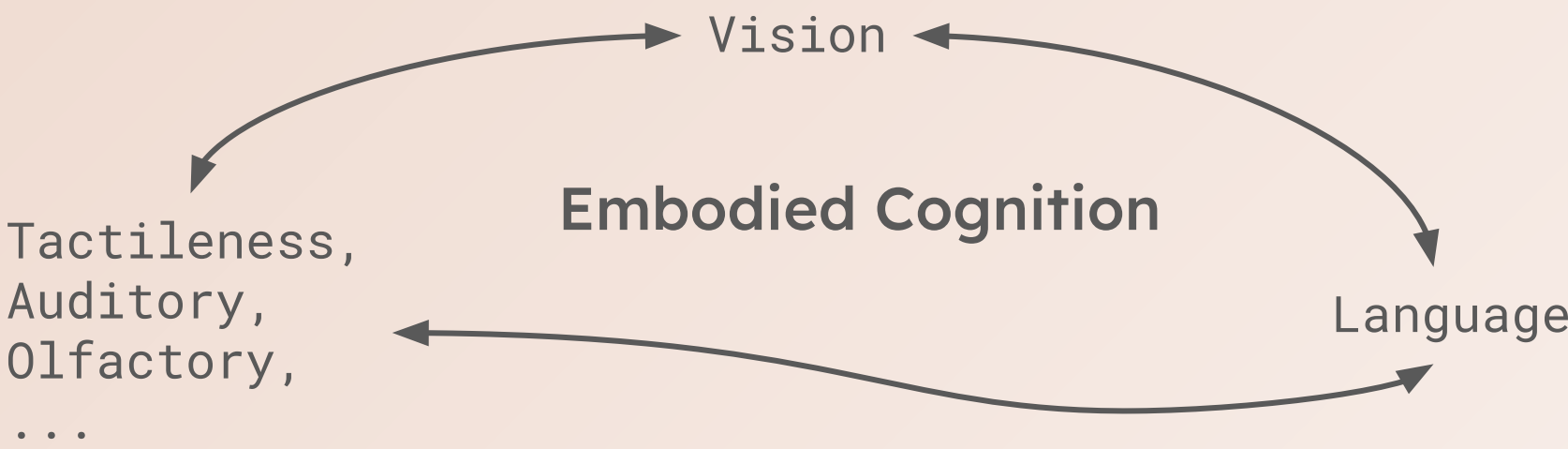
Lexical ambiguity



Semantic Clustering



Discussion / Future Work



Conclusion

- VG BERT still adheres to the classical NLP pipeline - 1.
- VG more strongly affects concrete concepts - 2.
- Better homonym distinction in VG embeddings - 3.
- VG embeddings better clustered into categories - 4.

References

[1] Zhang et al. Explainable Semantic Space by Grounding Language to Vision with Cross-Modal Contrastive Learning. 2021
[2] John R Searle et al. Minds, brains, and programs. The Turing Test: Verbal Behaviour as the Hallmark of Intelligence. 1980.
[3] J Ridley Stroop. Studies of interference in serial verbal reactions. 1935.
[4] Radford & Kim et al. Learning transferable visual models from natural language supervision. 2021
[5] Tenney et al. BERT Rediscovered the Classical NLP Pipeline. 2019