

UNIVERSITY OF  
CAMBRIDGE

# Measuring Cross-Modal Interactions in Multimodal Models

Laura Wenderoth<sup>1</sup>, Konstantin Hemker<sup>1</sup>, Nikola Simidjievski<sup>2,1</sup>, Mateja Jamnik<sup>1</sup>

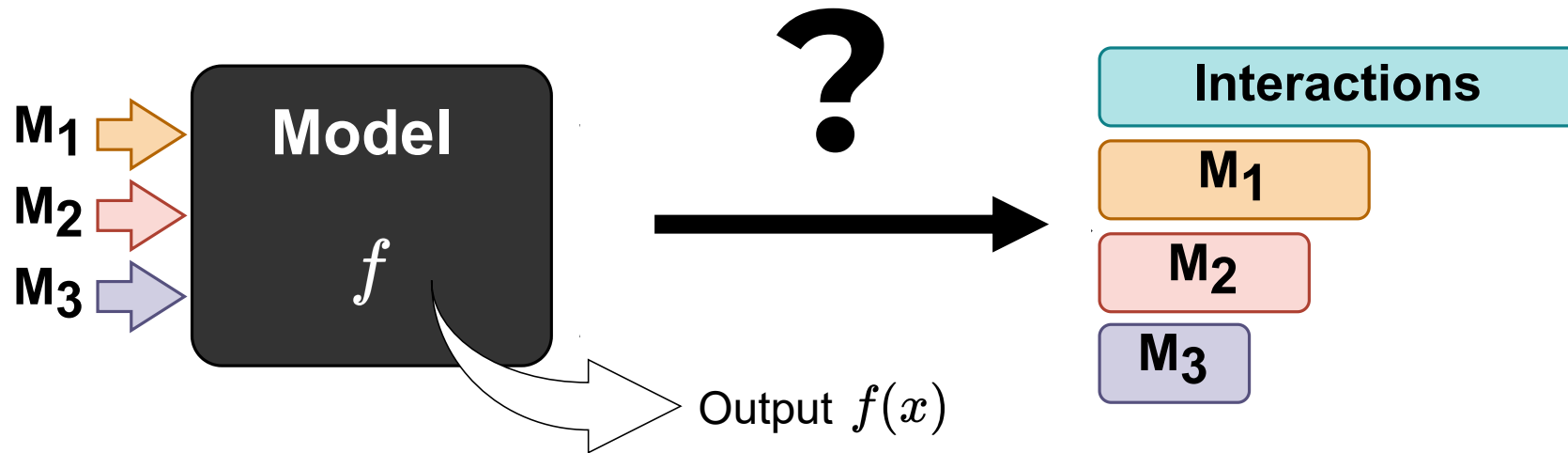
<sup>1</sup>Department of Computer Science and Technology, University of Cambridge, Cambridge, UK

<sup>2</sup>PBCI, Department of Oncology, University of Cambridge, Cambridge, United Kingdom

{lw457, kh701, ns779, mj201}@cam.ac.uk



# Motivation



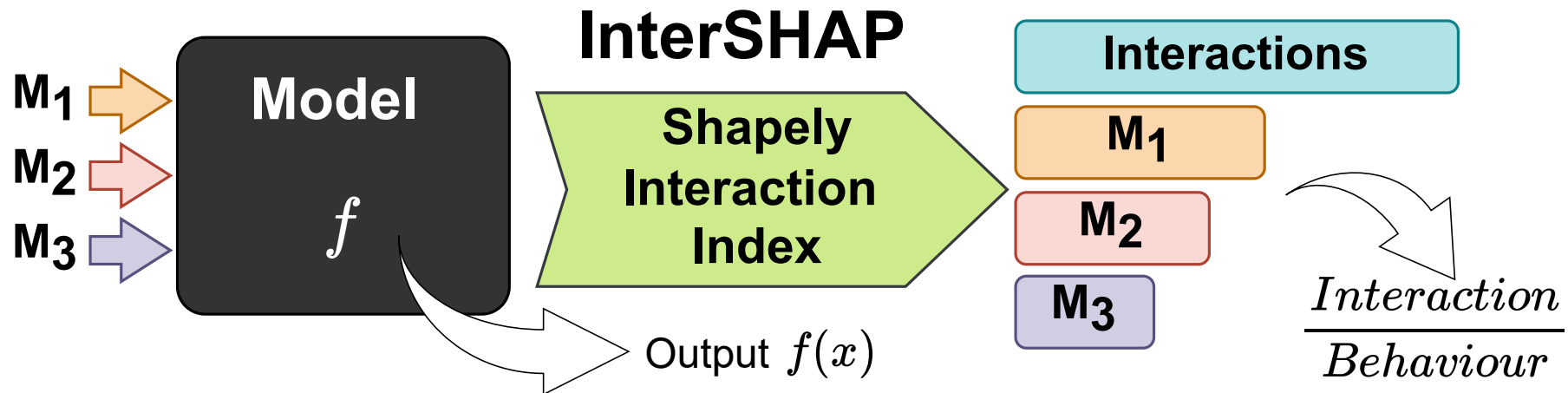
# Limitations of Previous Work

InterSHAP overcomes the limitations of other cross-modal interaction scores: it is unsupervised, performance agnostic, applicable to more than two modalities, and allows for dataset- (global) and sample-level (local) explainability

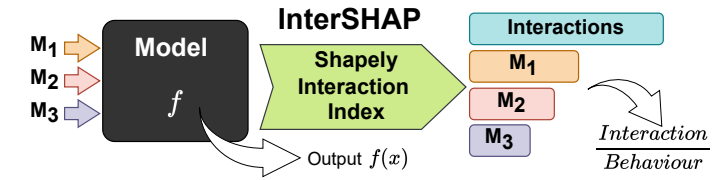
Score	Modalities > 2	Local	Unsupervised	Performance Agnostic
PID <sup>[1]</sup>	✗	✗	✓	✓
EMAP <sup>[2]</sup>	✗	✗	✗	✗
SHAPE <sup>[3]</sup>	✓	✗	✗	✗
<b>InterSHAP</b> (Ours)	✓	✓	✓	✓



# InterSHAP



# InterSHAP



$$\Phi_{ij} = \left| \frac{1}{N} \sum_{a=1}^N \phi_{ij}((m_1^a, \dots, m_M^a), f) \right|, \quad i, j \in \{1, \dots, M\} \quad (1)$$

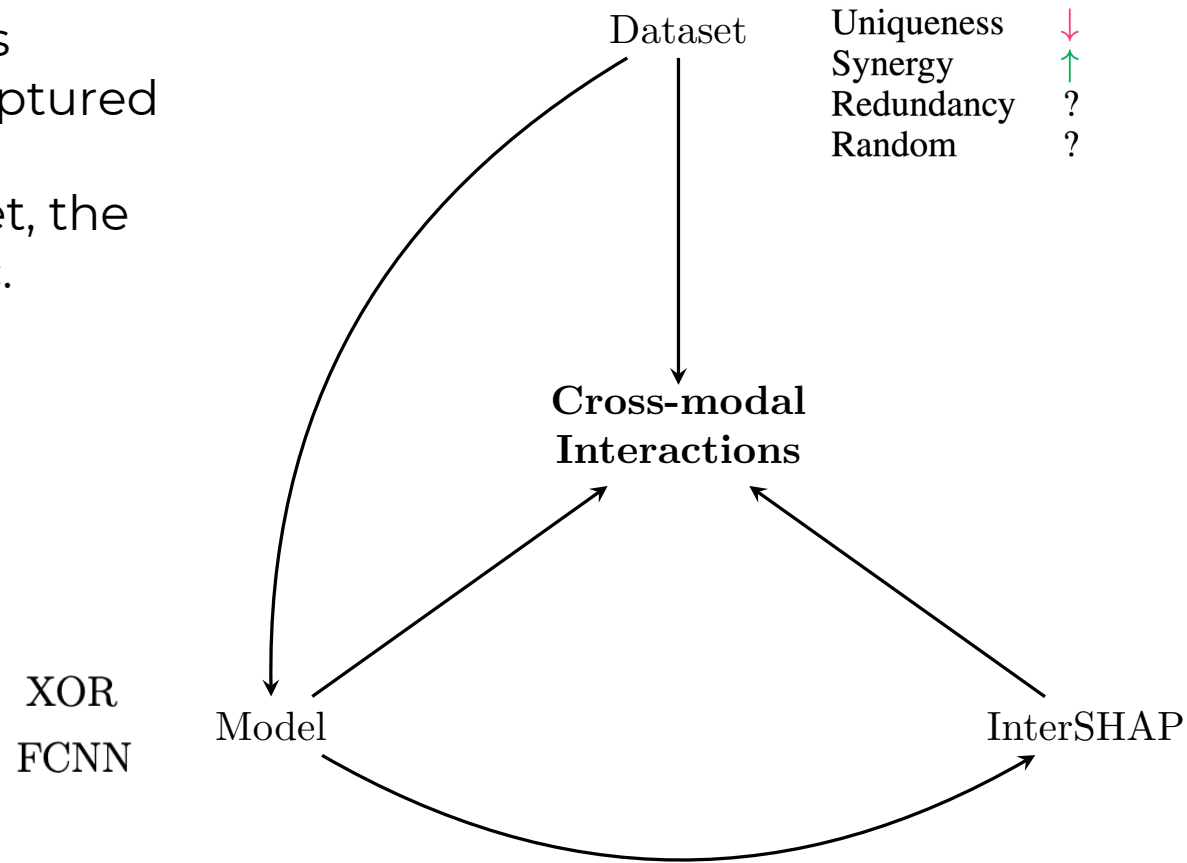
$$\Phi = \begin{bmatrix} \Phi_{11} & \Phi_{12} & \dots & \Phi_{1M} \\ \Phi_{21} & \Phi_{22} & \dots & \Phi_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_{M1} & \Phi_{M2} & \dots & \Phi_{MM} \end{bmatrix} \quad (2)$$

$$InterSHAP = \frac{Interactions}{Behaviour} \quad (3)$$



# Verification on Synthetic Data

Three main factors  
influencing the captured  
cross-modal  
interactionsdataset, the  
model, and metric.



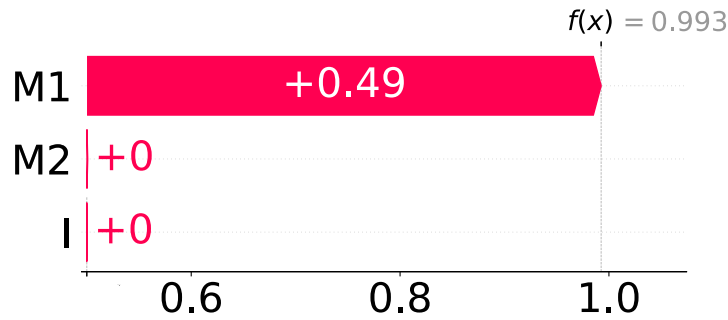
# Results – 2 Modalities

InterSHAP values are presented as percentages for both the XOR function and FCNN with early fusion on the HD-XOR datasets. Results for XOR align with expectations, confirming the effectiveness of InterSHAP. For the FCNN, slightly higher values for uniqueness and lower values for synergy suggest the FCNN model did not fully capture all underlying cross-modal interactions from the dataset.

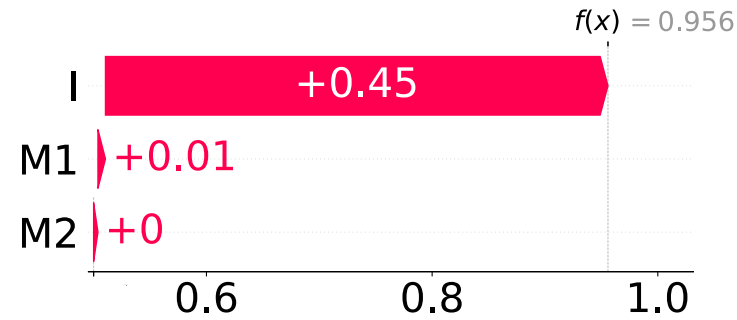
	Uniqueness		Synergy		Redundancy	Random
	XOR	FCNN	XOR	FCNN	FCNN	FCNN
InterSHAP	0.0	0.2 $\pm 0.1$	99.7	98.0 $\pm 0.5$	38.6 $\pm 0.5$	57.8 $\pm 1.1$



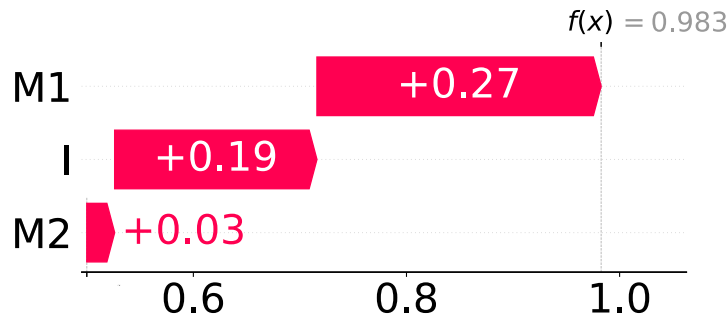
# Visualisation of Results



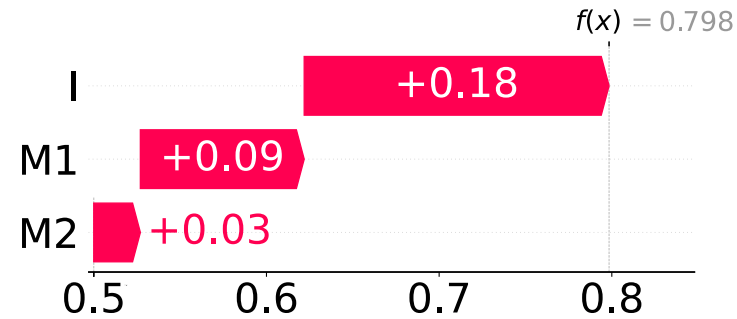
(a) Uniqueness



(b) Synergy



(c) Redundancy



(d) Random

Visualisation of results with FCNN with early fusion on the HD-XOR datasets. M1 represents modality 1, M2 modality 2, and I interactions...





# Results $> 2$ Modalities

	Uniqueness	Synergy	Redundancy
2 Modalities	0.2 $\pm 0.1$	98.0 $\pm 0.5$	38.6 $\pm 0.5$
3 Modalities	0.6 $\pm 0.2$	88.8 $\pm 0.5$	51.9 $\pm 0.3$
4 Modalities	1.2 $\pm 0.1$	64.1 $\pm 0.8$	40.2 $\pm 0.2$



# Limitations

- Runtime:  $O(N^M)$



# Application to healthcare domain

## Single Cell Dataset [4]

Modalities

- RNA
- Protein

Task: Cell Class Classification

Table 7: Details of Single Cell.

Class Distribution			
Neutrophil	Erythrocyte	B-Lymphocytes	Monocyte
43.7%	49.8%	1.4%	5.1%

## MIMIC III [5]

Modalities

- 12 physiological measurements (e.g. heart rate, 24h)
- static information on patients

Tasks: ICD and Mortality Classification

Table 8: Details of MIMIC III, ICD 1 and mortality.

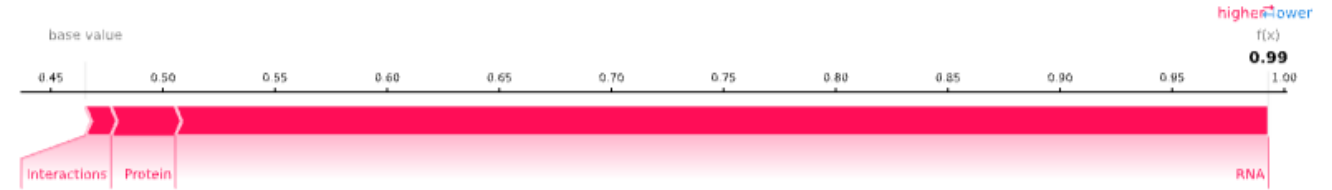
Class Distribution							
ICD		Mortality					
No	Yes	1d	2d	3d	7d	1 year	> 1 year
82.5%	17.5%	76.0%	0.4%	1.3%	1.0%	11.0%	10.3%



# Single Cell Dataset [4]

Table 4: Cross-modal interactions scores on the multimodal single-cell dataset for FCNN with early, intermediate and late fusion. InterSHAP aligns with other SOTA methods, capturing the decline in cross-modal information from early to late fusion.

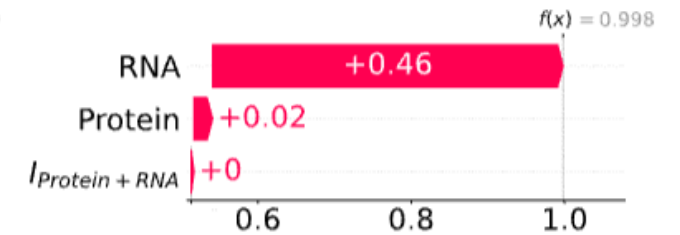
	Single-Cell		
	early	intermediate	late
<b>InterSHAP</b>	$1.9 \pm 0.4$	$1.5 \pm 0.4$	$0.4 \pm 0.1$
PID	$0.08 \pm 0.01$	$0.08 \pm 0.01$	$0.06 \pm 0.0$
EMAP <sub>gab</sub>	$0 \pm 0$	$0 \pm 0$	$0 \pm 0$
SHAPE	$1.0 \pm 0.2$	$0.7 \pm 0.2$	$0 \pm 0$



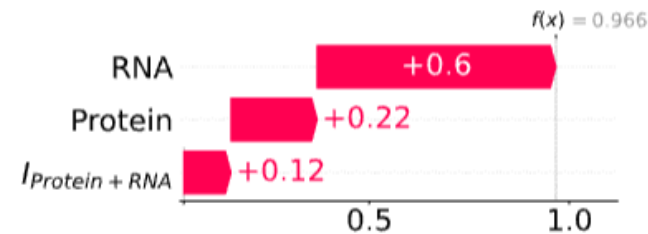
(a) Explanation over the whole dataset



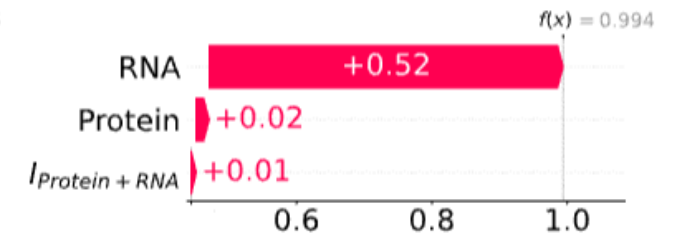
(b) B-Cell Progenitor



(c) Erythrocyte Progenitor



(d) Monocyte Progenitor



(e) Neutrophil Progenitor



# MIMIC III [5]

Table 5: Cross-modal interactions scores on the MIMIC III dataset for baseline model and MVAE model from Multi-Bench implementation. InterSHAP aligns with other SOTA methods, capturing greater cross-modal interaction in the baseline compared to MVAE, while uniquely quantifying the proportional contribution of cross-modal interactions.

	ICD-9		Mortality	
	baseline	MVAE	baseline	MVAE
<b>InterSHAP</b>	1.2 $\pm$ 0.2	6.8 $\pm$ 1.3	11.0 $\pm$ 0.5	12.3 $\pm$ 2.8
PID	0.06 $\pm$ 0.01	0.09 $\pm$ 0.01	0.10 $\pm$ 0.01	0.11 $\pm$ 0.01
EMAP <sub>gap</sub>	0 $\pm$ 0	1.2 $\pm$ 0.0	-0.8 $\pm$ 0.1	0.9 $\pm$ 0.1
SHAPE	0.2 $\pm$ 0	0.6 $\pm$ 0	0.2 $\pm$ 0.2	0.7 $\pm$ 0.2



# Summary

- Novel cross-modal interaction score: **InterSHAP** - Open-Source implementation with integration into SHAP package
  - >2 Modalities
  - Local
  - Unsupervised
  - Performance agnostic
- Application to healthcare multimodal datasets

More in our paper:

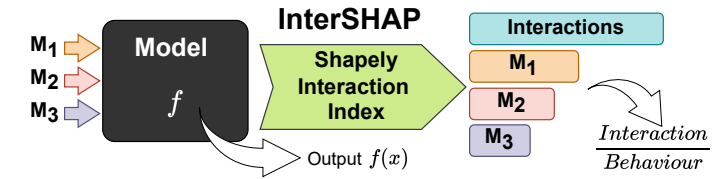
- Quantitative evaluation of existing cross-modal interaction scores



# References

- [1] Hessel, J.; and Lee, L. 2020. Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think! In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 861–877. Association for Computational Linguistics.
- [2] Hu, P.; Li, X.; and Zhou, Y. 2022. SHAPE: An Unified Approach to Evaluate the Contribution and Cooperation of Individual Modalities. In Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, 3064– 3070. Vienna, Austria: International Joint Conferences on Artificial Intelligence Organization.
- [3] Liang, P. P.; Cheng, Y.; Fan, X.; Ling, C. K.; Nie, S.; Chen, R.; Deng, Z.; Allen, N.; Auerbach, R.; Mahmood, F.; Salakhutdinov, R. R.; and Morency, L.-P. 2023. Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. Advances in Neural Information Processing Systems, 36: 27351–27393.
- [4] Burkhardt, D.; Luecken, M.; Benz, A.; Holderrieth, P.; Bloom, J.; Lance, C.; Chow, A.; and Holbrook, R. 2022. Open Problems - Multimodal Single-Cell Integration.
- [5] Liang, P. P.; Lyu, Y.; Fan, X.; Wu, Z.; Cheng, Y.; Wu, J.; Chen, L.; Wu, P.; Lee, M. A.; Zhu, Y.; Salakhutdinov, R.; and Morency, L.-P. 2021. MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. In Vanschoren, J.; and Yeung, S.-K., eds., Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021.

# Shapley Interaction Index



$$\phi_{ij}(M, f) = \sum_{S \subseteq M \setminus \{i, j\}} \frac{|S|!(M - |S| - 2)!}{2(M - 1)!} \nabla_{ij}(S, f), \quad i \neq j.$$

$$\nabla_{ij}(S, f) = [f_{S \cup \{ij\}}(S \cup \{ij\}) - f_{S \cup \{i\}}(S \cup \{i\}) - f_{S \cup \{j\}}(S \cup \{j\}) + f_S(S)]$$

$$\phi_{ii}(M, f) = \phi_i(M, f) - \sum_{j \in M} \phi_{ij}(M, f) \quad \forall i \neq j.$$

$$\phi_i(M, f) = \sum_{S \subseteq M \setminus \{i\}} \frac{|S|!(|M| - |S| - 1)!}{|M|!} \Delta$$

$$\Delta = [f_{S \cup \{i\}}(S \cup \{i\}) - f_S(S)].$$