

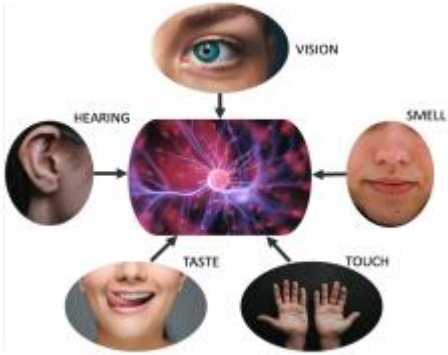
Data representation for multimodal models

Lidia Pivovarova

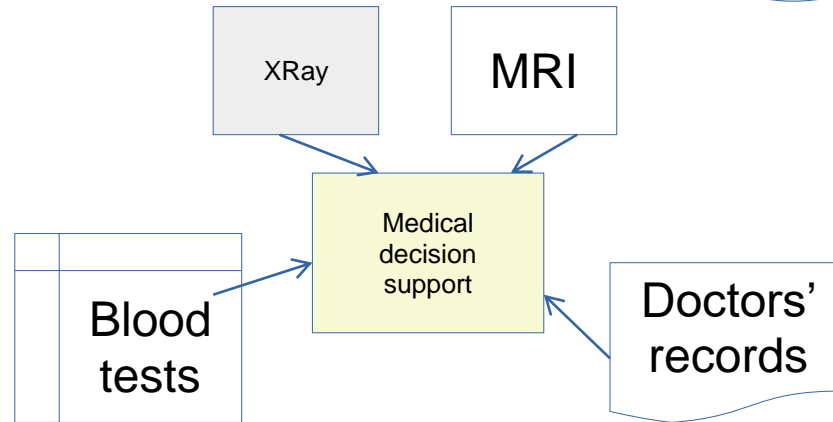
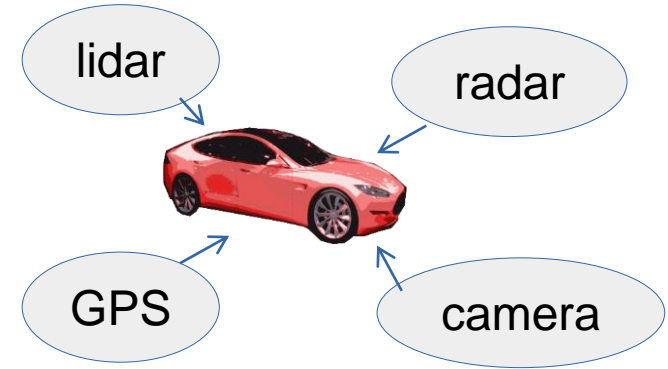
28.08.2024

Summer study group

Multimodality in machine learning



(Zhao et al., 2024)



Application domains

- .Robotics
- .Healthcare
- .Multimedia
- .Human-computer interaction
- .Interactive agents
- .etc

Principles of multimodality

1. Modalities are heterogenous

- diverse in quality, structure, distributions, noise, and relevance

2. Modalities are connected

- statistical and semantic relations, redundancy and uniqueness of information

3. Modalities interact

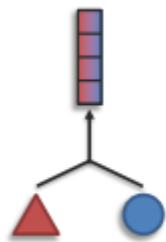
modalities processed together result in a novel

Unimodal representations

- Shift from hand-crafted features to data-driven, obtained as a result of self-supervised pretraining
 - Visual features: convolutional networks, e.g. ResNet,
 - Textual features: embeddings, LLMs
 - Audio: neural acoustic models
 - Graphs: either linearization and LLMs or graph networks

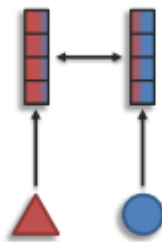
Representation

Fusion



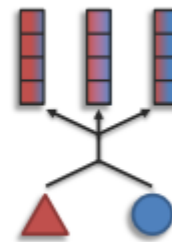
modalities > # representations

Coordination



modalities = # representations

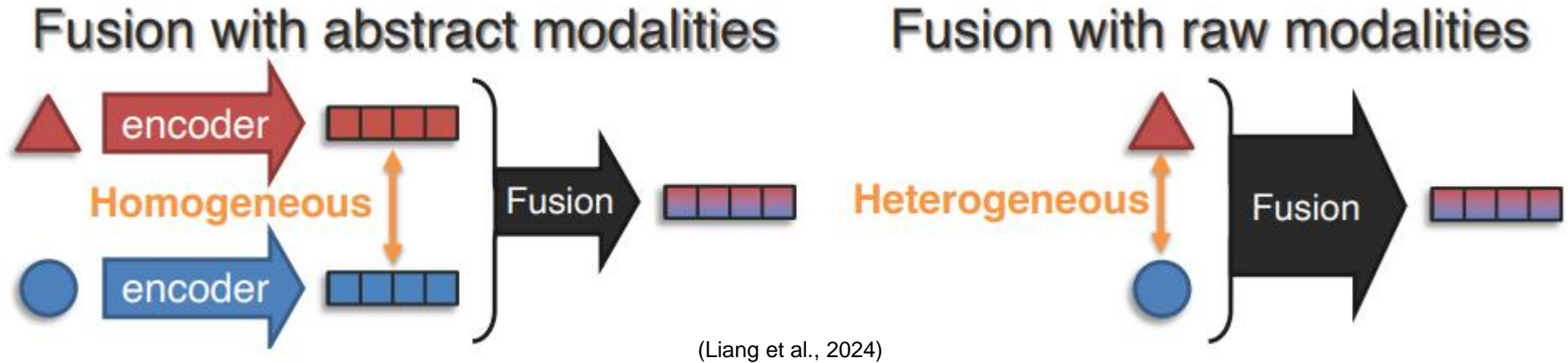
Fission



modalities < # representations

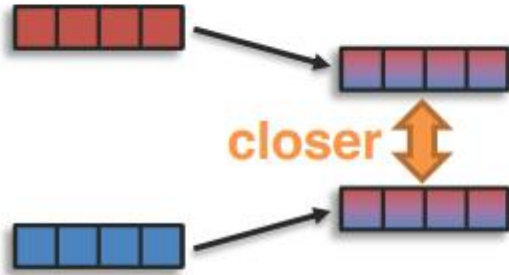
(Liang et al., 2024)

Fusion

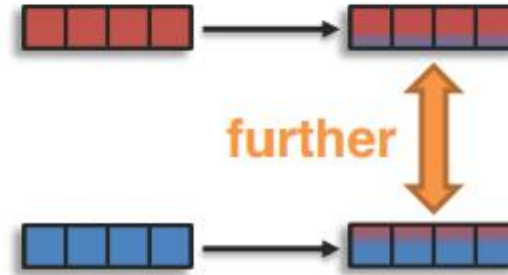


Coordination

Strong coordination



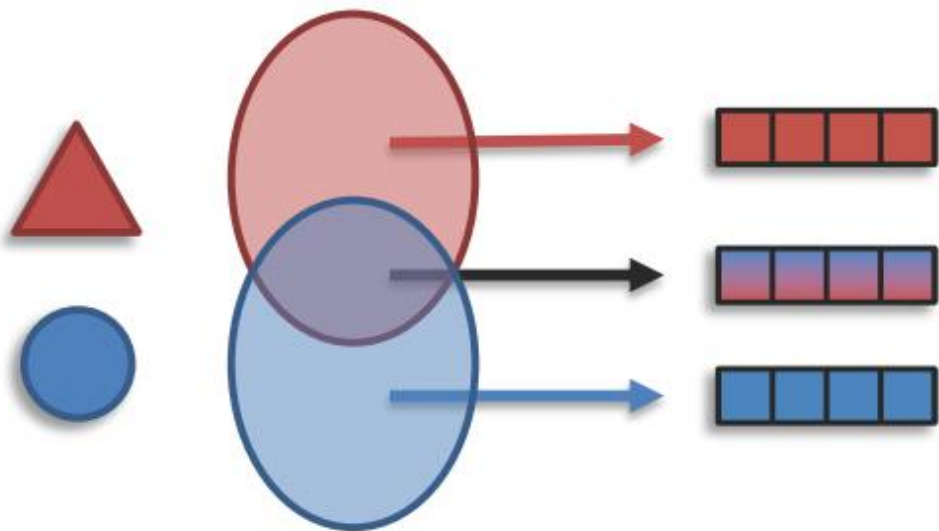
Partial coordination



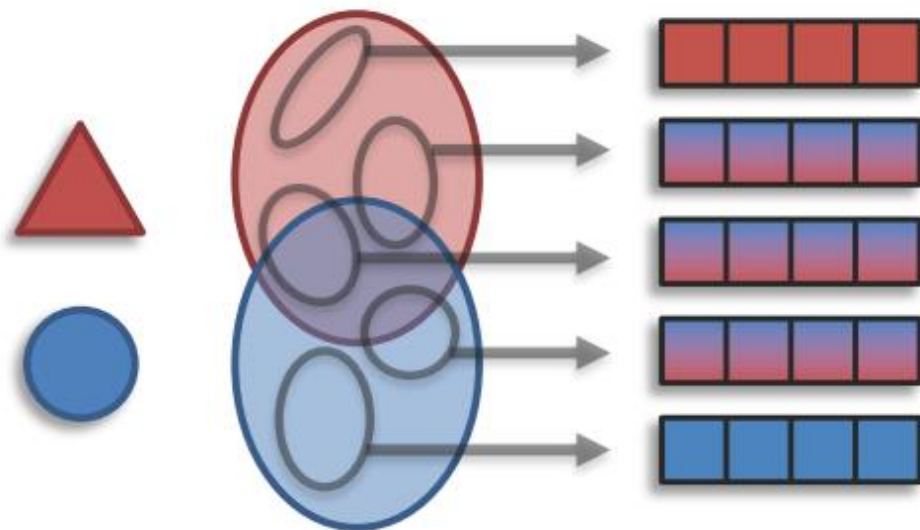
(Liang et al., 2024)

Fission

Modality-level fission



Fine-grained fission



(Liang et al., 2024)

Typical tasks

Frameworks	Applications	Modalities
Joint representation	Video classification	A + V
		A + V + T
	Event detection	A + V
		A + V + T
	Sentiment analysis	A + V + T
	Visual question answering	I + T
	Emotion recognition	A + V
A + V + T		
Speech recognition	A + V	
Coordinated representation	Cross-modal retrieval	I ~ T
	Image caption	I ~ T
	Cross-modal embeddings	V ~ T
		I ~ T
Transfer learning	I ~ T	
Encoder-decoder	Image caption	I → T
	Video description	V → T
	Text to image synthesis	T → I

(Guo et al., 2019)

Document understanding



T H E

CHARTER of FORESTS (a)

Granted by King JOHN to his
Subjects in the Year 1215*.



JOHN, by the Grace of God, King of
England, &c. Know ye, that for the
Honour of God, and the Health
of our Soul, and the Souls of our
Ancestors and Successors, and for the
Exaltation of Holy-Church, and for the
Reformation of our Kingdom, We
have of our free and good Will given and granted for
Us and our Heirs, these Liberties hereafter specified, to
be had and observ'd in our Kingdom of England for
ever.

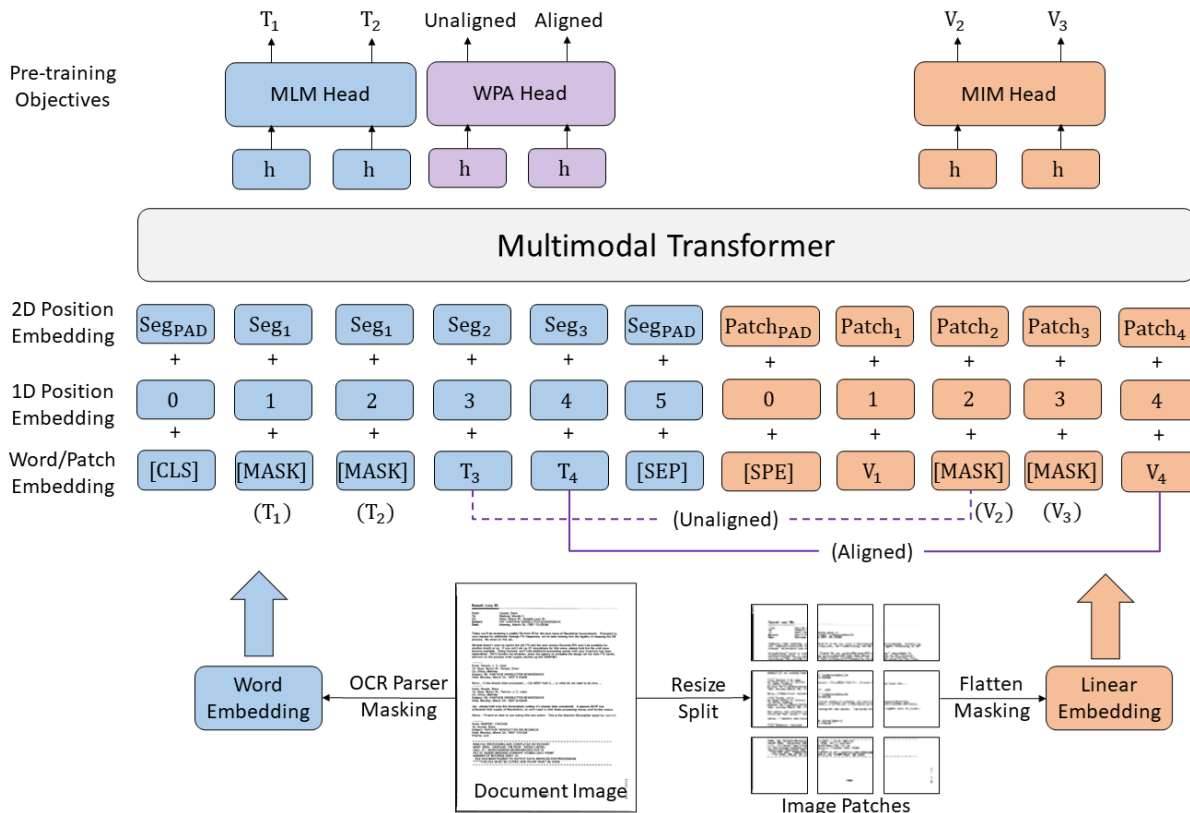
VOL. III.

Mm 2

I. Im.

(a) The Forests belong'd originally to the Crown, and the Kings had granted several Parts and Parcels to private Men, who had grubb'd them up and made them Arable, or Pasture. But yet all that was thus grubb'd was still call'd Forest. These Forests belonging to the King as his own Demesnes, or as the Sovereign Lord, were a continual Source of vexatious Suits, as well against those which held them of the King, as against the neighbouring Freemen under pretence of the Rights of the Crown.

* As it is to be found in Matthew Paris, p. 250.



Sources

- Liang, P. P., Zadeh, A., & Morency, L. P. (2024). Foundations & trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Computing Surveys*, 56(10), 1-42.
- Zhao, F., Zhang, C., & Geng, B. (2024). Deep Multimodal Data Fusion. *ACM Computing Surveys*, 56(9), 1-36.
- Guo, W., Wang, J., & Wang, S. (2019). Deep