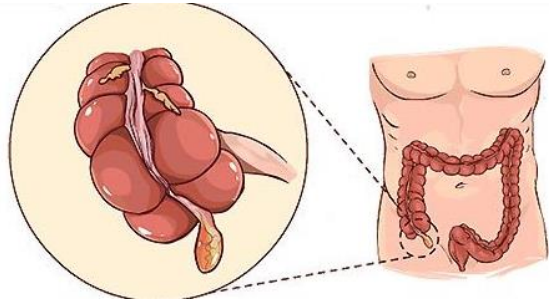

Multimodal Learning for Appendicitis Diagnosis

Saroj Baral

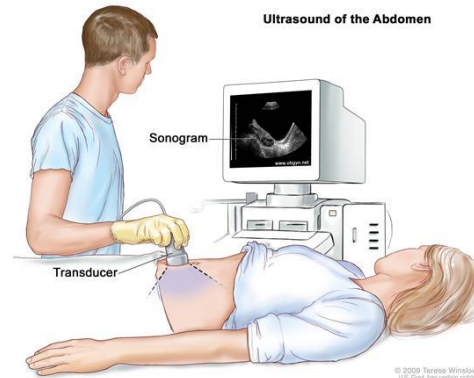
April 2, 2025



Introduction



Borrowed from indushealthplus (a)



Borrowed from nibib (b)



Borrowed from mayoclinic(c)

- **Appendicitis:** Inflammation of Appendix
- Generally, occurs in young individuals (ages 10–30)

Ultrasound

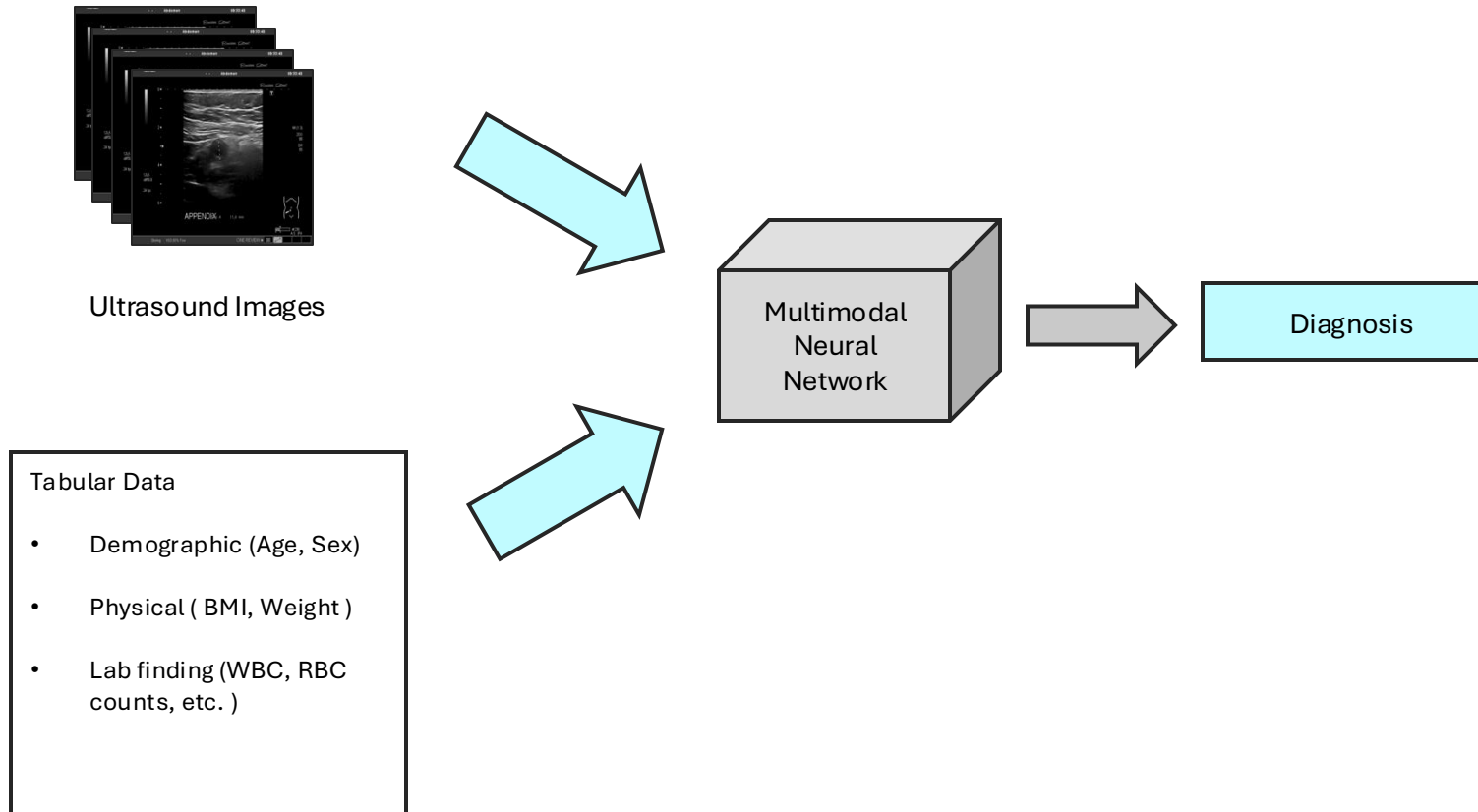
- Less expensive
- Uses Sound Wave
- No Radiation Exposure

Computed Tomography

- More expensive than Ultrasound
- Uses X rays to make detail images
- Radiation Exposure

a) https://www.indushealthplus.com/front/media/article_img/appendicitis-causes-symptoms-prevention.jpg
b) https://www.nibib.nih.gov/sites/default/files/inline-images/Ultrasound_Terese%20Winslow.jpg
c) <https://www.mayoclinic.org/tests-procedures/ct-scan/about/pac-20393675>

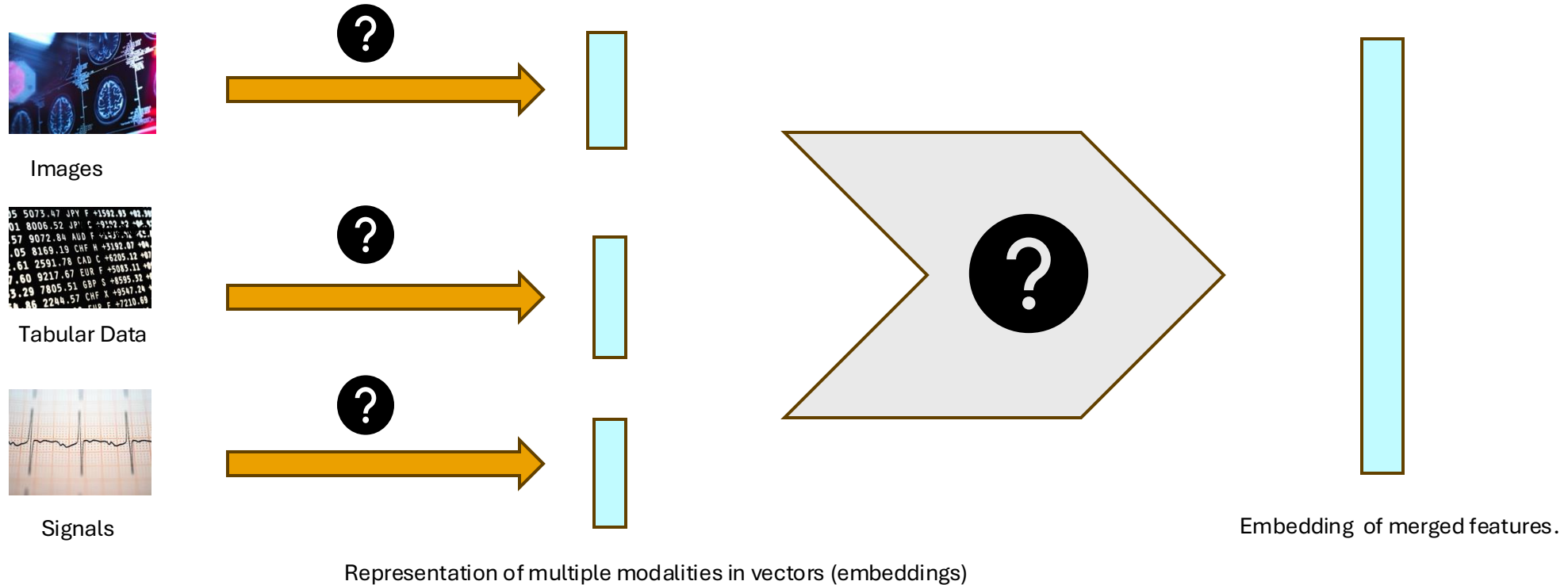
Overall Goal and Challenges



Major Challenges

- Image preprocessing
- Small data size (579 records)
- Varying no of images (views) per record

Challenges in Multimodal Architecture



Representation

How to encode data from different modality?

Fusion

How to combine information from different modality?

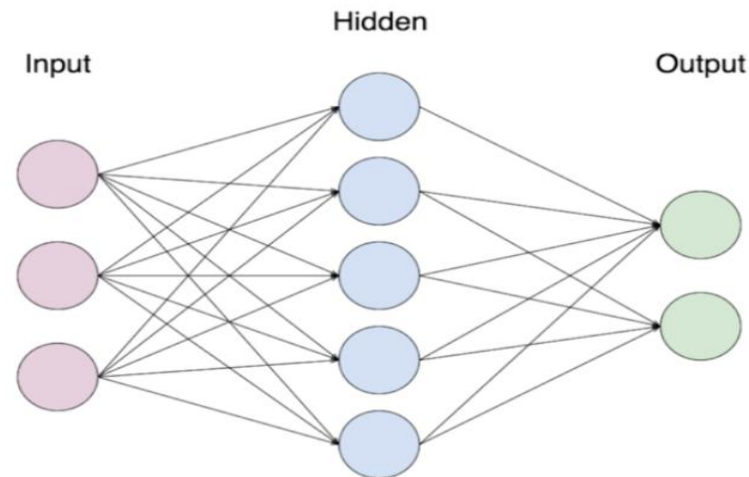
Proposed Model Architecture : Representing Tabular data

Input

- Lab reports (WBC, RBC count ,etc.)
- Physical (BMI, Weight)

Output

- Vector representing above feature (embedding)



Borrowed from specbee (d)

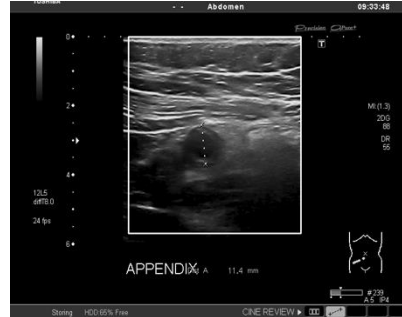
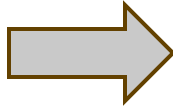
Neural Network extracts useful information from data

Image preprocessing



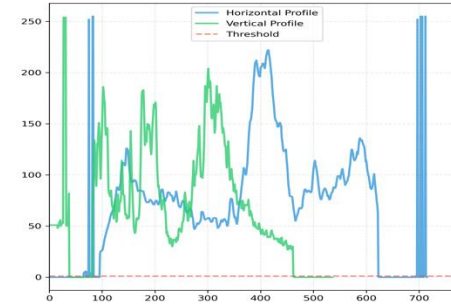
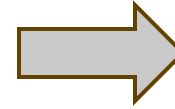
1. Original Image

- Annotations from ultrasound machine
- Markers by doctors
- Multiple in one file



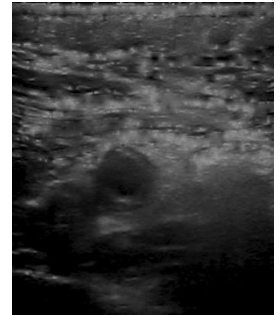
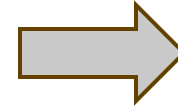
2. Original Image with boundaries

- Tried automated boundary detection
- Didn't work well for all
- Manually inspected and corrected



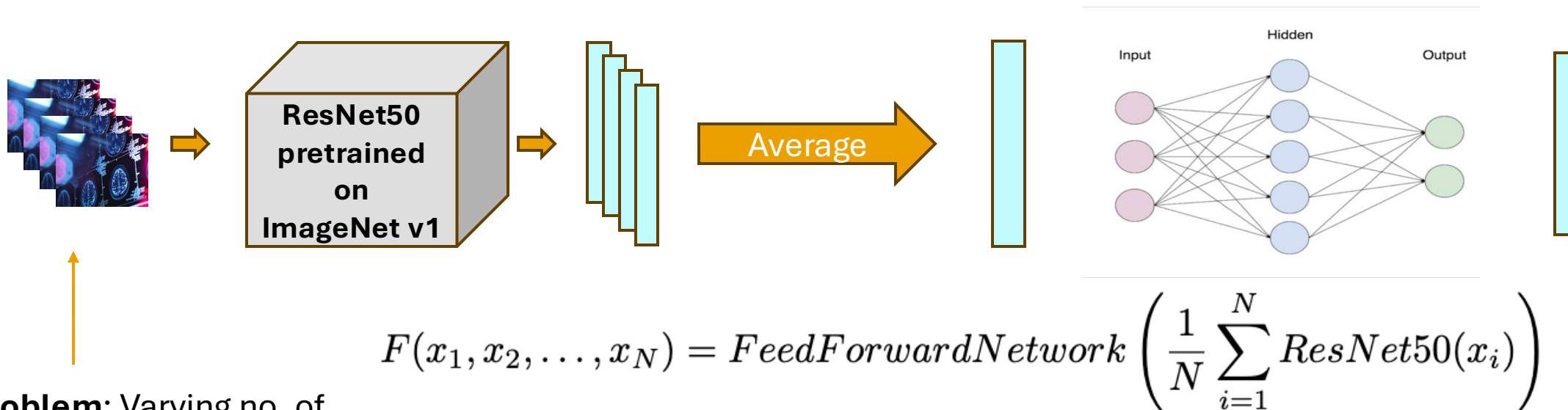
3. Intensity Profile

- Intensity profile is inspected
- Created mask for appropriate threshold



3. Final Image with annotations removed

Proposed Model Architecture : Representing Images



Problem: Varying no. of pictures of a single person

From Deep Sets by Zaheer et.al [e]

Theorem 2 A function $f(X)$ operating on a set X having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in X , iff it can be decomposed in the form $\rho \left(\sum_{x \in X} \phi(x) \right)$, for suitable transformations ϕ and ρ .

Model Architecture : Fusion Strategy

Fusion is way of combining the modality after representation.

Two main strategies:

- Early Fusion (Feature level fusion)
 - Fuse at input to the model
- Late fusion
 - Fuse later in intermediate or end (ensemble)

Late fusion are slightly effective. [f]

Using Late fusion in the proposed model



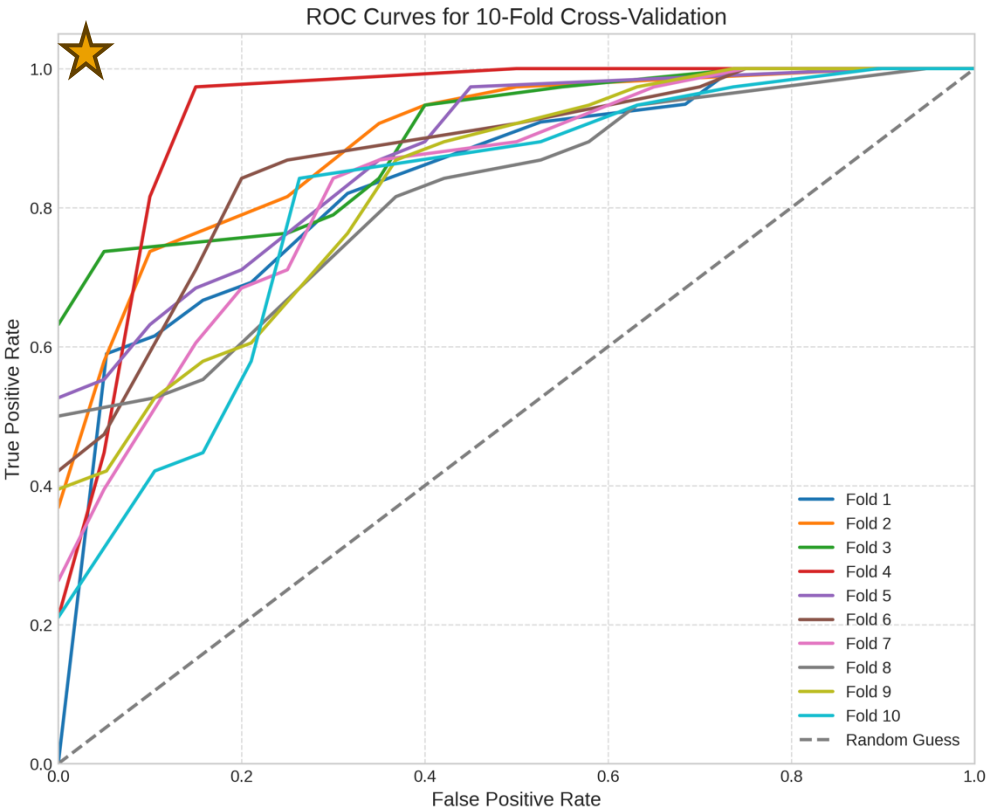
1. Simple Concatation

2. Weighted Sum (For interpretability)

Fused vector is then fed into another neural network for classification !!

Results

Perfect Model



Comparing with related work (**AUROC**) on same dataset

Multimodal	Images only (Marcinkevicius et. al [g])
0.84 ± 0.04	0.80 ± 0.06

Fusion Weights in 2nd fold

- $W1 = 1.045$ (tabular)
- $W2 = 0.034$ (images)

Conclusions



This multimodal neural network shows promising potential for improving diagnostic accuracy over single-modal methods.



Doctors are better at distinguishing ultrasound images than neural networks.



The approach can be extended to other medical conditions that require combining imaging and tabular data.



Future work can explore alternative fusion techniques or models and larger datasets for further improvements



**Thank you
for listening**
