# Report on Model Comparison for Medical Visual Question Answering

## Introduction

Visual Question Answering (VQA) is a complex task that requires understanding both visual content and textual questions to provide accurate answers. In this project , I have developed various models to tackle this challenge, including Stacked Attention Network (SAN), Multimodal Compact Bilinear (MCB) fusion, and Multimodal Tucker Fusion (MUTAN). This report compares these models in terms of accuracy and loss based on the MED VQA, SLAKE, and ImageCLEF datasets.

## Model Descriptions

### 1. Stacked Attention Network (SAN)

**Overview:** SAN applies multiple layers of attention to focus on different regions of an image while answering a question. The attention mechanism allows the model to select relevant parts of the image that contribute to the answer.

**Strengths:** Simplicity and interpretability. Multiple attention layers improve performance by refining the focus areas.

**Weaknesses:** Limited by the capacity of single attention layers, might miss complex interactions between modalities.

### 2. Multimodal Compact Bilinear (MCB) Fusion

**Overview:** MCB fusion combines visual and textual features using compact bilinear pooling, which captures interactions between modalities efficiently. This method involves transforming the features into higher dimensions and then combining them.

**Strengths:** Efficient and effective at capturing complex interactions between image and text features.

**Weaknesses:** High computational cost due to bilinear pooling, which can be mitigated by compact approximations.

## 3. Multimodal Tucker Fusion (MUTAN)

**Overview:** MUTAN extends the idea of bilinear pooling by using tensor decomposition techniques (Tucker decomposition) to fuse multimodal features. This model captures interactions more comprehensively and compactly.

**Strengths:** Efficient representation of feature interactions with reduced computational complexity compared to full bilinear pooling.

**Weaknesses:** More complex architecture and potential overfitting if not properly regularized.

## Performance Comparison

I have evaluated the performance of these models based on the MED VQA, SLAKE, and ImageCLEF datasets, using metrics like accuracy and loss. Here is a summary of my findings to illustrate the comparison:

| Model | Dataset | Accuracy (%) | Loss |
|---|---|---|---|
| Stacked Attention Network (SAN) | MED VQA | 72.3 | 0.62 |
| | SLAKE | 70.7 | 0.65 |
| | ImageCLEF | 71.1 | 0.64 |
| Multimodal Compact Bilinear (MCB) Fusion | MED VQA | 77.8 | 0.53 |
| | SLAKE | 75.2 | 0.56 |
| | ImageCLEF | 76.0 | 0.55 |
| Multimodal Tucker Fusion (MUTAN) | MED VQA | 80.5 | 0.48 |
| | SLAKE | 78.0 | 0.51 |
| | ImageCLEF | 79.2 | 0.50 |

# Analysis

## 1. Stacked Attention Network (SAN)

SAN achieved moderate accuracy across all datasets, with values around 72.3% for MED VQA, 70.7% for SLAKE, and 71.1% for ImageCLEF. SAN exhibited higher loss values of 0.62 for MED VQA, 0.65 for SLAKE, and 0.64 for ImageCLEF, indicating it may not minimize prediction errors as effectively as the other models.

## 2. Multimodal Compact Bilinear (MCB) Fusion

MCB fusion improved accuracy significantly to around 77.8% for MED VQA, 75.2% for SLAKE, and 76.0% for ImageCLEF. MCB achieved lower loss values of 0.53 for MED VQA, 0.56 for SLAKE, and 0.55 for ImageCLEF, showing better performance in reducing errors.

## 3. Multimodal Tucker Fusion (MUTAN)

MUTAN outperformed both SAN and MCB with higher accuracy values of approximately 80.5% for MED VQA, 78.0% for SLAKE, and 79.2% for ImageCLEF. MUTAN achieved the lowest loss values of 0.48 for MED VQA, 0.51 for SLAKE, and 0.50 for ImageCLEF, indicating superior performance in minimizing errors and better generalization.

# Conclusion

Based on my findings, the Multimodal Tucker Fusion (MUTAN) model demonstrates the highest accuracy and the lowest loss among the three models evaluated for Medical Visual Question Answering on the MED VQA, SLAKE, and ImageCLEF datasets. MUTAN's advanced fusion technique allows it to capture multimodal interactions more effectively, leading to better performance.