

# "Решение задачи VQA посредством нейросетевой модели"

Дипломная работа Базовая кафедра распознавания изображений и  
обработки текста АВВУУ

Студент: Даниял Алиев  
Руководитель: Артур Бегаев

# Соревнование Receipt-VQA-2023

## Входные данные:

- Изображение чека
- Вопрос
- Категория вопроса ("amount", "count", "ratio")
- Валюта, используемая в чеке
- Список операций необходимый для вычисления результата

## Выходные данные:

Ответы на вопросы в формате float

## Метрика:

$$\text{MASE} = \text{mean} \left( \frac{|e_j|}{\frac{1}{J} \sum_{j=1}^J |Y_j - \bar{Y}|} \right) = \frac{\frac{1}{J} \sum_j |e_j|}{\frac{1}{J} \sum_j |Y_j - \bar{Y}|}$$

# Соревнование Receipt-VQA-2023

Input image



Input question:

How much should be paid?

receipt\_currency: Malaysian ringgit

question\_category: amount

operations: amount

Output:  
6.2

# Соревнование Receipt-VQA-2023

Input image

SYARIKAT PERNIAGAAN GIN KEE  
(81109-A)  
NO 290, JALAN AIR PANAS,  
SETAPAK,  
53200, KUALA LUMPUR,  
TEL 03-40210276  
GST ID : 000750673920  
SIMPLIFIED TAX INVOICE

CASH  
Doc No: CS00012524 Date: 03/01/2018  
Cashier: USER Time: 17:08:00  
Salesperson: Ref:

| Item                           | Qty | S/Price | Amount | Tax |
|--------------------------------|-----|---------|--------|-----|
| 1512                           | 1   | 10.60   | 10.60  | SR  |
| 104 COTTON GLOVE (DOZEN)       |     |         |        |     |
| 3032                           | 1   | 23.32   | 23.32  | SR  |
| FACE MASK                      |     |         |        |     |
| 3313                           | 2   | 19.08   | 38.16  | SR  |
| CS 200A CUTTING WHEEL          |     |         |        |     |
| Total Qty                      | 4   |         | 72.08  |     |
| Total Sales (Excluding GST)    |     |         | 68.00  |     |
| Discount                       |     |         | 0.00   |     |
| Total GST                      |     |         | 4.08   |     |
| Rounding                       |     |         | 0.00   |     |
| Total Sales (Inclusive of GST) |     |         | 72.08  |     |
| CASH                           |     |         | 72.08  |     |
| Change                         |     |         | 0.00   |     |

GST SUMMARY

| Tax Code | % | Amnt (RM) | Tax (RM) |
|----------|---|-----------|----------|
| SR       | 6 | 68.00     | 4.08     |
| Total    |   | 68.00     | 4.08     |

GOODS SOLD ARE NOT RETURNABLE. THANK YOU.

Input answer:

How much is tax to total amount ratio?,

Receipt\_currency: Malaysian ringgit

question\_category: ratio

Operations: division

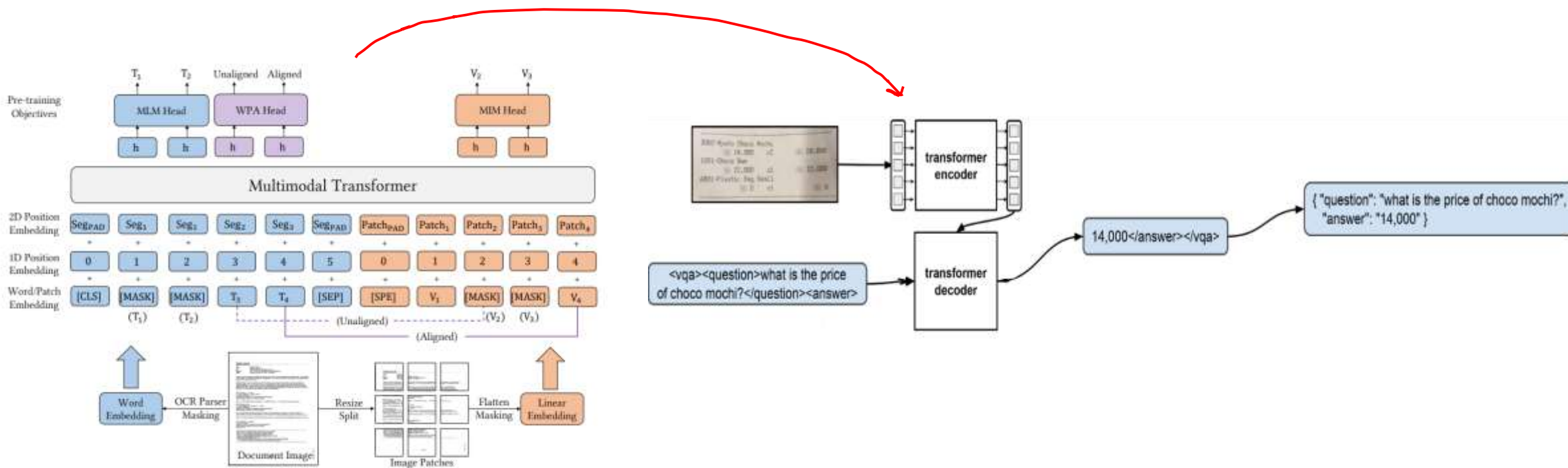
Output:

$$4.08/72.08=0.0566$$

# Работа модели на категории “amount”

Категория amount включает в себя вопросы, ответы на которые содержатся в самом изображении – не требуется никаких дополнительных вопросов.

В качестве модели для fine-tuning берем модель **LayoutLMv3** ([LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking](#))



# Работа модели на категории “ratio”

Ratio определяется как отношение двух чисел, каждое из которых можно найти, используя уже обученную LayoutLMv3 модель

What is **net total** and **total amount** ratio?

lmV3

lmV3

|  |        |            |           |
|--|--------|------------|-----------|
| RESTORAN WAN SHENG<br>002043319-W<br>No.2, Jalan Temenggung 19/9,<br>Seksyen 9, Bandar Mahkota Cheras,<br>43200 Cheras, Selangor<br>GST REG NO: 001335787520 |        |            |           |
| Tax Invoice  |        |            |           |
| INV No.: 1057981 Cashier: Thandar<br>Date : 23-03-2018 13:28:55  |        |            |           |
| Description  | Qty    | U.price    | Total TAX |
| Teh (B)  | 1 x    | 2.20       | 2.20 SR   |
| Cham (B)   | 1 x    | 2.20       | 2.20 SR   |
| Bunga Kekwa  | 1 x    | 1.70       | 1.70 SR   |
| Take Away  | 3 x    | 0.20       | 0.60 SR   |
| Total QTY: 6   |        |            |           |
| Total (Excluding GST):   |        |            | 6.33      |
| GST payable (6%):  |        |            | 0.37      |
| Total (Inclusive of GST):  |        |            | 6.70      |
| TOTAL :  |        |            | 6.70      |
| CASH :   |        |            | 6.70      |
| GST Summary  |        |            |           |
| SR   | (# 6%) | Amount(RM) | Tax(RM)   |
|  |        | 6.33       | 0.37      |

output

$$6.33/6.9 = 0.944$$

# Работа модели на категории “count”

Категория достаточно объемная – рассмотрим для начала вопросы, которые подразумевают некоторые арифметические операции. Тогда алгоритм будет следующим:

- 1) Определить шаблон выражения, которое даст ответ на вопрос (с помощью operations)
- 2) Находим неизвестные с помощью предобученной LayoutLMv3
- 3) Подставляем числа в выражение и получаем результат

How much should be paid for 3 products from the 2nd position?

Receipt\_currency: Malaysian ringgit  
question\_category: count  
Operations: multiplication; sorting

| Item                | Qty | Price | Amount | Tax |
|---------------------|-----|-------|--------|-----|
| 2576 HOES PIN       | 3   | 1.05  | 3.15   | SR  |
| 1937 4' HOES HANDLE | 1   | 7.95  | 7.95   | SR  |
| Total Qty           | 4   |       | 11.13  |     |

Total Sales (Exclusive of GST) 10.50  
Discount 0.00  
Total GST 0.63  
Rounding 0.00  
Total Sales (Inclusive of GST) 11.13  
CASH 11.13  
Change 0.00

GST SUMMARY  
Tax Code % Amt (RM) Tax (RM)  
SR 0 10.50 0.63  
Total 10.50 0.63

$X * Y$

$7.95 * 3 = 23.85$

lmV3

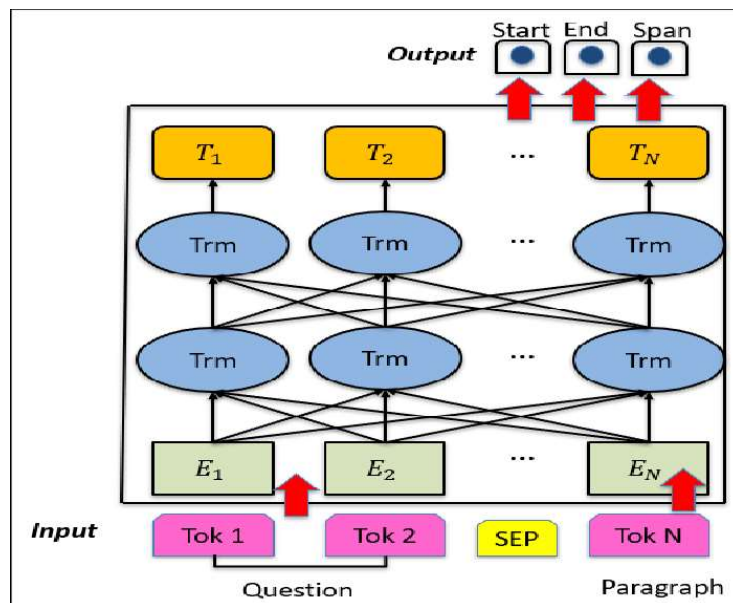
# Извлечение информации из текста

Для ответа на вопросы, требующие арифметические операции, мы извлекали ключевые слова из вопроса и использовали их для подачи в LayoutLmv3.

Для выполнения данной задачи нужно решить задачу question answering – для этого нужно предварительно подготовить ключевые слова для каждого из вопросов, используя регулярные выражения (все вопросы одной категории очень похожи)

Далее снова воспользуемся трансформером с Hugging Face – BertForQuestionAnswering - для fine-tuning на наших данных.

How many goods in the 1st position can be purchased for 22?



1<sup>st</sup> position



22

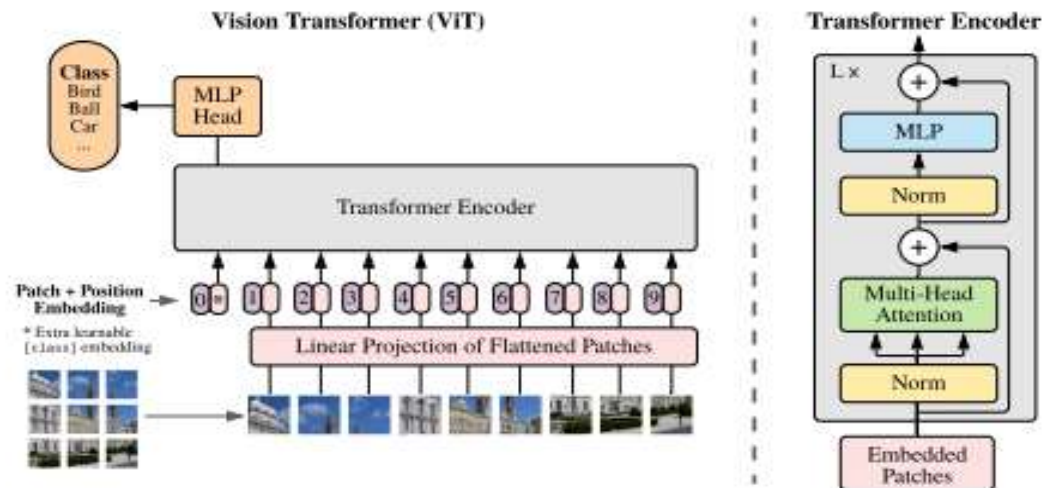


# Работа модели на категории “count”

В данной категории присутствуют так же вопросы, ответы на которые не получаются вычислением арифметических выражений, пример: **How many positions are free of cost?**

Ответы на такие вопросы дать достаточно тяжело – по итогу лучший вариант оказалось использование ViTModel ([https://huggingface.co/docs/transformers/model\\_doc/vit#transformers.ViTModel](https://huggingface.co/docs/transformers/model_doc/vit#transformers.ViTModel))

Используем last\_hidden\_state, который подается на линейный слой и Softmax – MLP Head (используем классификацию на 8 элементов – на train выборке этот диапазон покрывается 90% ответов)



# Результаты

## VQA Track

| Participant       | Best Submission Name | MASE Total | MASE Amount | MASE Count | MASE Ratio |
|-------------------|----------------------|------------|-------------|------------|------------|
| daniyallaiev      | submission0          | 0.7874     | 0.7979      | 0.5718     | 0.9926     |
| BASLINE (abegaev) | answer               | 0.8786     | 0.8068      | 0.8291     | 1.0000     |

Анализ полученных результатов:

- На amount результаты практически идентичны, потому что baseline модель так же использовала **LayoutLMv3** на этой категории
- Видим существенное улучшение на категории count за счет грамотного применения Bert и предобученной на прошлой категории модели
- На ratio результат так же превосходит baseline

Проблема модели:

- На многих изображениях датасета OCR работает некорректно и зачастую “не видит” ответ на вопрос, после чего дальнейшие действия уже бесполезны

Возможные решения:

- Самостоятельно обучение OCR на исходном датасете (требуется дополнительной разметки)
- Предварительная обработка изображений – применение фильтров и т.д.

# Дальнейшие действия

- Улучшение работы модели – корректировка работы OCR, применение предварительной обработки изображений, применение новых идей архитектуры модели
- Проверить работу модели на **TAT-DQA** датасете (<https://github.com/NExTplusplus/TAT-DQA>)
- Проанализировать полученные результаты, найти уязвимости модели
- Снова постараться улучшить модель

**Спасибо за внимание!**