# NTU DLCV (Autumn, 2022) HW3 Report

R10921069　沈郁鈞

## Problem 1

### 1

Methods analysis (3%)

The main reason is that CLIP changes the task from **image classification** to **text-image pairing**.
Thus, CLIP has higher performance than other *image classification based* models on **unseen class** task.
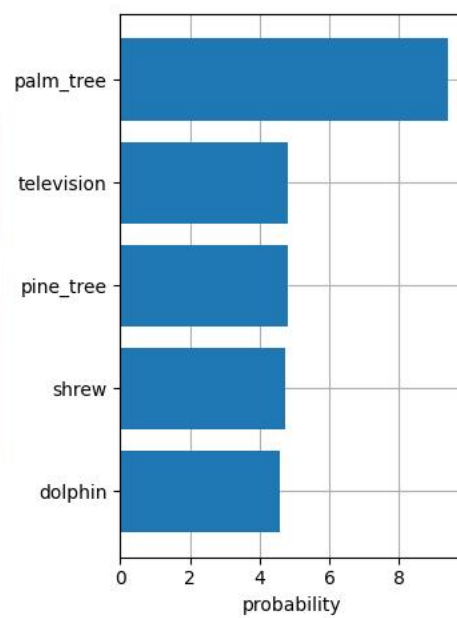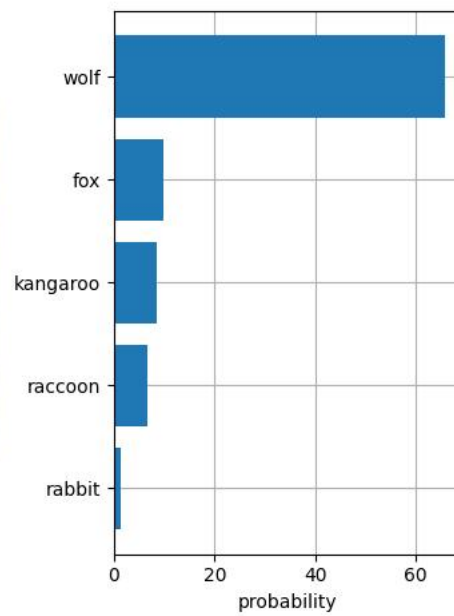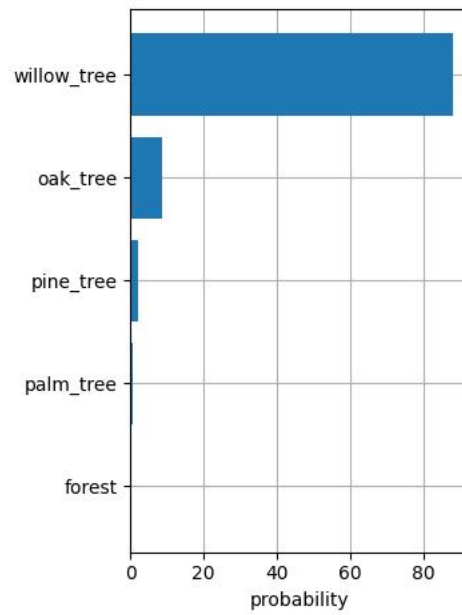
### 2

| Prompt templates | Accuracy | Rank |
|---|---|---|
| This is a photo of {object} | 0.6876 | 1 |
| This is a {object} image. | 0.6824 | 2 |
| No {object}, no score. | 0.5644 | 3 |

#### Discussion

I think the prompt template should be as simple and accurate as we could. Also, the template with positive sentence pattern would improve the performance.

### 3

## Problem 2

# 1

Decoding Strategy: Greedy (k=1)
learning rate: $10^{-5}$
epoch: 6

| CIDEr | CLIPScore |
|---|---|
| 0.959487 | 0.731331 |

# 2

## SETTING 1

Decoding Strategy: Greedy (k=1)
learning rate: $2*10^{-5}$
epoch: 6

| CIDEr | CLIPScore |
|---|---|
| 0.86034 | 0.73010 |

## SETTING 2

Decoding Strategy: beam search (k=3)
learning rate: $2*10^{-5}$
epoch: 6

| CIDEr | CLIPScore |
|---|---|
| 0.87933 | 0.71645 |

## SETTING 3

Decoding Strategy: Greedy (k=1)
learning rate: $10^{-5}$
epoch: 3

| CIDEr | CLIPScore |
|---|---|
| 0.87437 | 0.72860 |

# Problem 3

# 1

`BIKE.PNG`

| [BOS] | A | person | riding | a |
|---|---|---|---|---|
| | | | | |

| bicycle | on | a | city | street |
|---|---|---|---|---|
| | | | | |

| . | [EOS] |
|---|---|
| | |

**GIRL.PNG**

| [BOS] | A | little | girl | is |
|---|---|---|---|---|
| | | | | |

| eating | a | slice | of | pizza |
|---|---|---|---|---|
| | | | | |

| . | [EOS] |
|---|---|
| | |

**SHEEP.PNG**

Two sheep standing in a field of grass .

**SKI.PNG**



A man and a woman standing in the snow with a ski poles .

**UNBRELLA.PNG**

A woman and a woman are standing under an umbrella. [EOS]

## 2

**TOP-1**

CLIPScore:0.977172



A woman holding a kite in a field. [EOS]

**LAST-1**

CLIPScore:0.427246



| [BOS] | A | man | wearing | a |
| --- | --- | --- | --- | --- |
| hat | and | a | hat | is |
| standing | in | a | large | room |
| . | [EOS] | | | |

## 3

### TOP-1 PREDICTION

- The predicted caption is reasonable for the picture.
  - (The kite is really on the woman's hands!!!)
- I found that **"kite"** is almost precisely marked on the position in the picture.

### LAST-1 PREDICTION

- The predicted caption is partially correct and partially unreasonable.
  - **O** It's true that the man is wearing a helmet(hat) on his head.
  - **X** It's impossible for a hat to stand.
  - **?** I don't know the exact place where the man stands in.