



Class-Token



Language-Token



Learnable Weight



Image Embedding



Frozen Parameters



Learnable Parameters











Input

CLIP Vision Encoder

Query
 q $\text{sim}(k, q)$
Eq. (3)

CLIP Text Encoder

  + <class>
  + <class>
  + <class>
 ...
  + <class>

First-level Prompt

 k

Class Prototypes

 W_1
 W_2
 W_3
 \vdots
 W_n
 \otimes

Second-level Prompt

Select Top- K
Eq. (4)

Selected Prompt

Regularization
Eq. (5)

Selected Prompt

Prompt Selection

Prediction

Semantic KD
Eq. (8)

Knowledge Distillation

Multi-Head
Self attention
Layers

VIT

Input

CLIP Text Encoder

First-level Prompt

Semantic
Feature
Space