# **CountGD**: Multi-Modal Open-World Counting
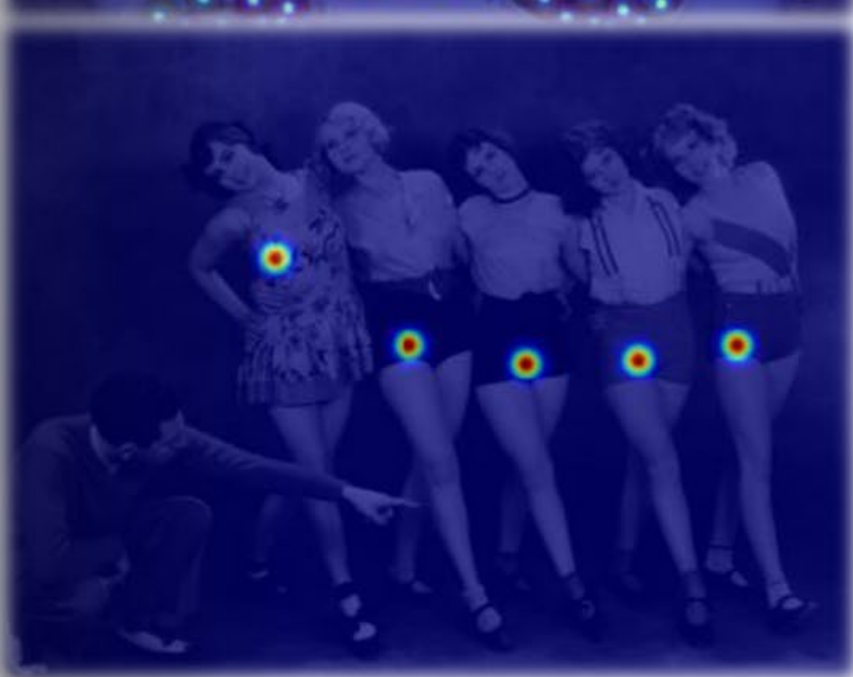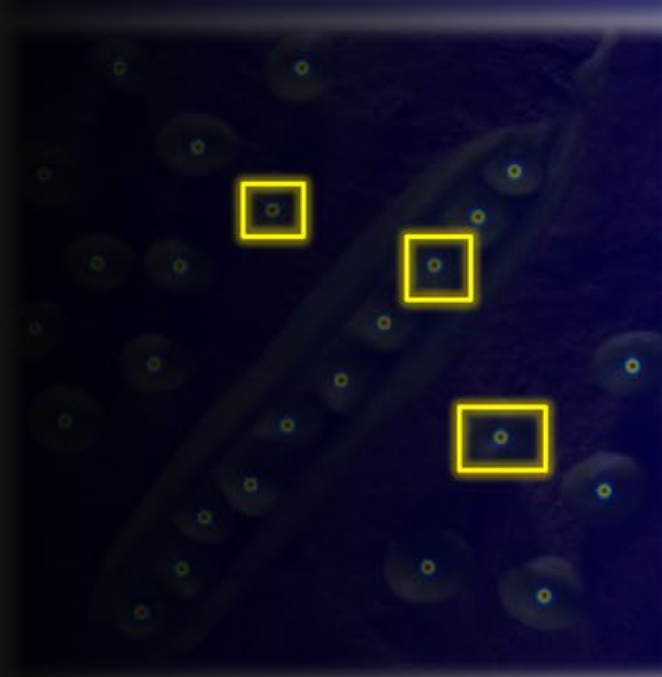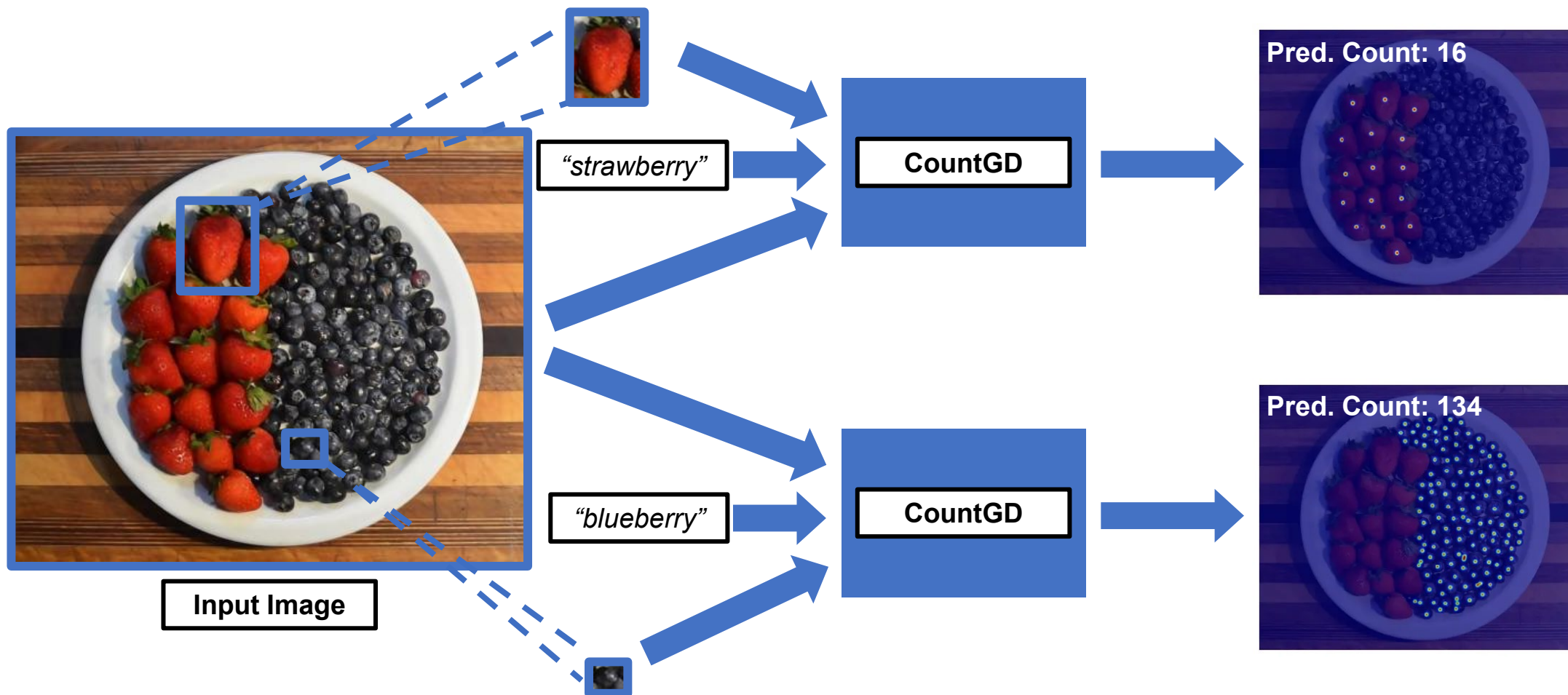
## NeurIPS 2024

*Niki Amini-Naieni, Tengda Han, &
Andrew Zisserman*

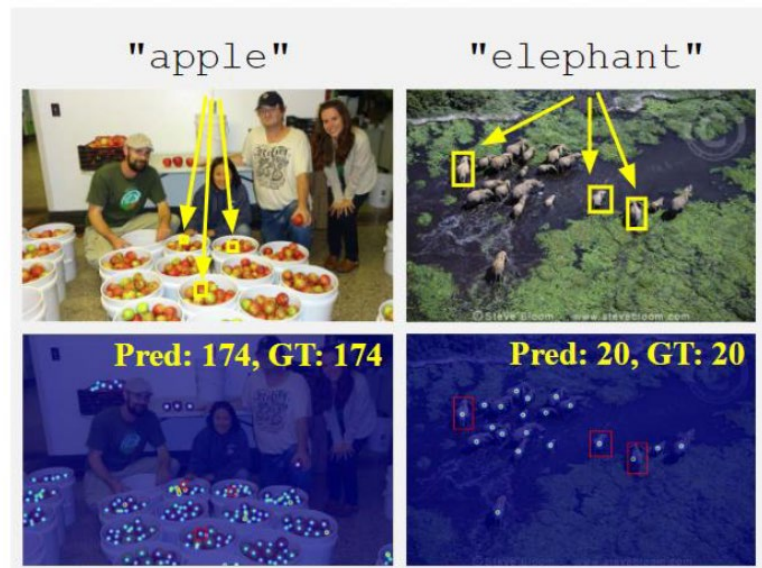[Project Page](#)| [ArXiv](#) | [Code](#) | [App](#)

# Main Idea

# Contributions

1. We introduce the first open-world counting model, CountGD, where the prompt can be specified by a text description or visual exemplars or both.
2. We show the performance of CountGD significantly improves the state-of-the-art on multiple counting benchmarks.
3. We carry out a preliminary study into different interactions between the text and visual exemplar prompts, including the cases where they reinforce each other and where one restricts the other.



**Visual Exemplars & Text**
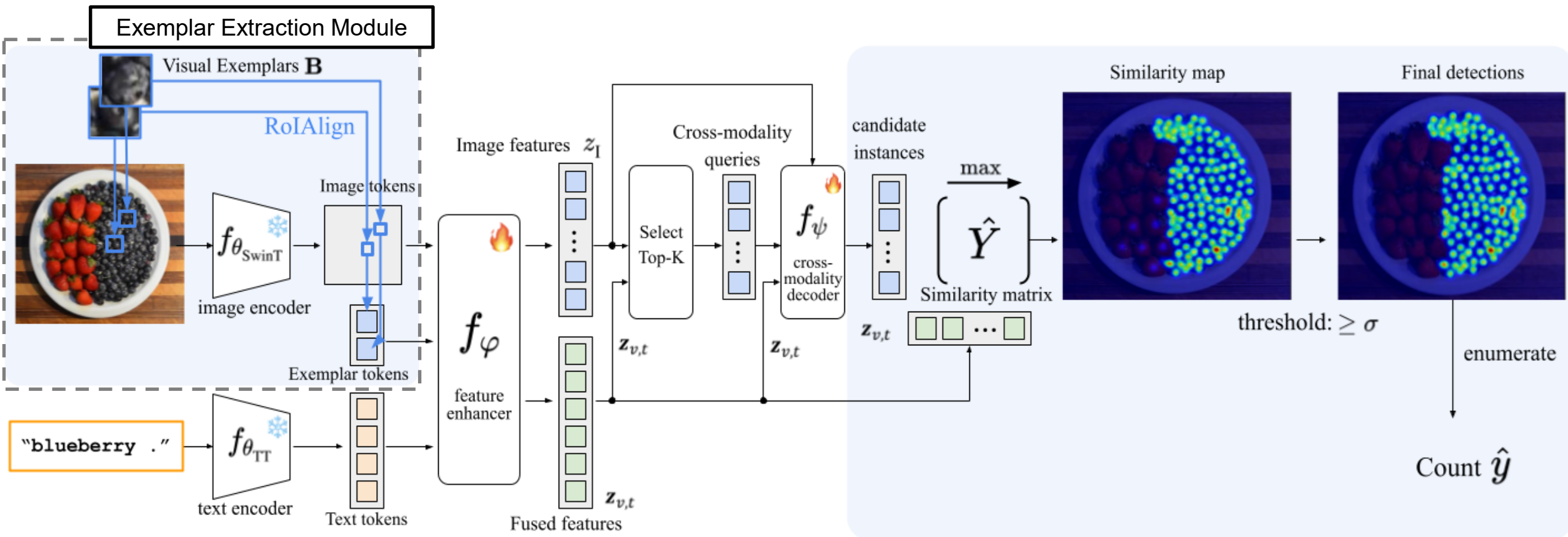
Input text is in quotes, and input visual exemplars are boxed

**Text Only**

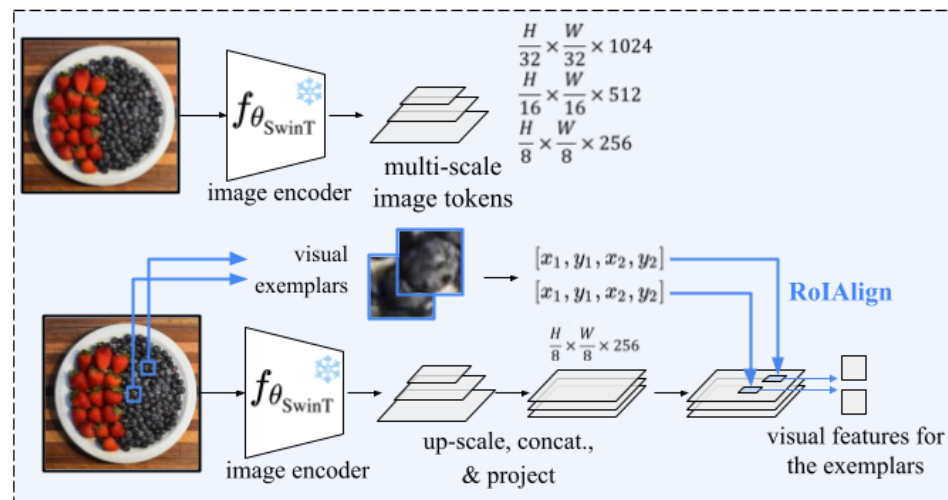...tions are plotted and overlaid on top of eac...

# CountGD Architecture

: Added to GroundingDINO
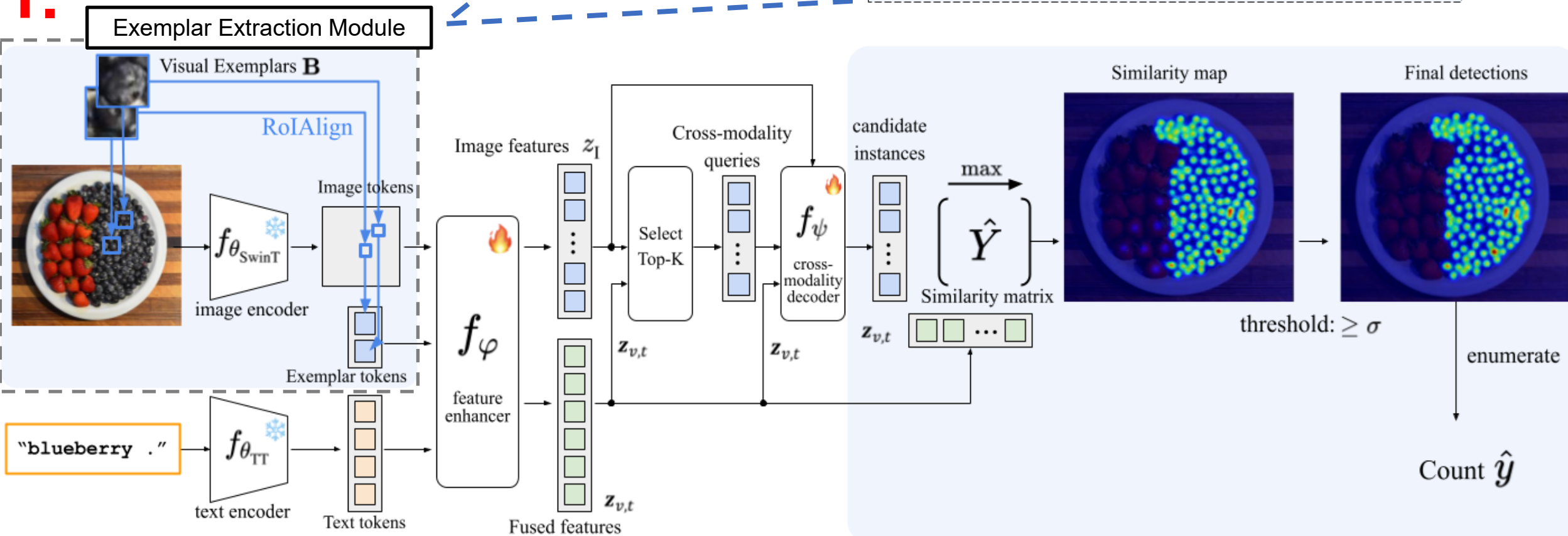
Module to enable inputting visual exemplars into GroundingDINO.

1.

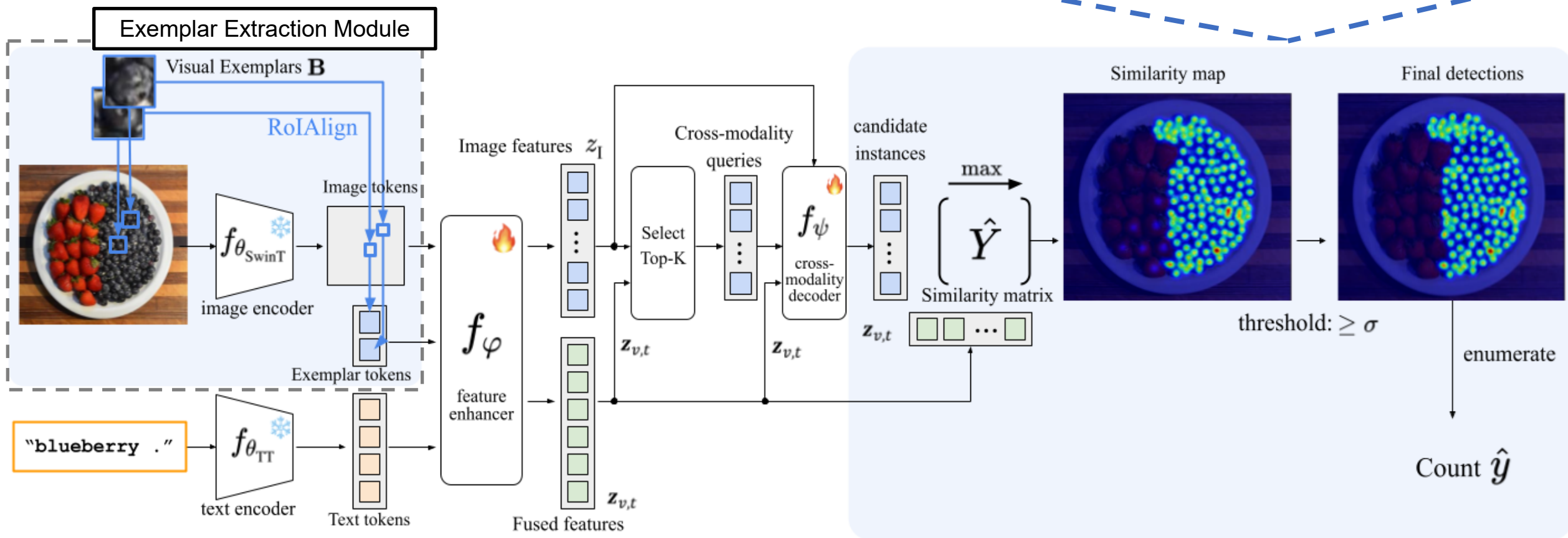$\frac{H}{32} \times \frac{W}{32} \times 1024$
$\frac{H}{16} \times \frac{W}{16} \times 512$
$\frac{H}{8} \times \frac{W}{8} \times 256$

image encoder $f_{\theta_{SwinT}}$ — multi-scale image tokens

visual exemplars → $[x_1, y_1, x_2, y_2]$ → $[x_1, y_1, x_2, y_2]$ → RoIAlign

image encoder $f_{\theta_{SwinT}}$ → up-scale, concat., & project $\frac{H}{8} \times \frac{W}{8} \times 256$ → visual features for the exemplars

1.

Exemplar Extraction Module

Visual Exemplars $\mathbf{B}$

RoIAlign

image encoder $f_{\theta_{SwinT}}$ → Image tokens → Exemplar tokens

"blueberry ." → text encoder $f_{\theta_{TT}}$ → Text tokens

feature enhancer $f_{\varphi}$ → Fused features $z_{v,t}$

Image features $z_I$ → Select Top-K → Cross-modality queries → cross-modality decoder $f_{\psi}$ → candidate instances

$z_{v,t}$ → $z_{v,t}$ → $z_{v,t}$

Similarity matrix → $\hat{Y}$ → max → Similarity map → threshold: $\geq \sigma$ → Final detections → enumerate → Count $\hat{y}$

# Getting the final count.

**2.**



## Exemplar Extraction Module

Visual Exemplars $\mathbf{B}$

RoIAlign

Image tokens

$f_{\theta_{SwinT}}$ ❄️

image encoder

Exemplar tokens

Image features $z_I$

Cross-modality queries

candidate instances

Similarity map

Final detections

Select Top-K

$f_\psi$ 🔥 cross-modality decoder

$\xrightarrow{\max}$

$\begin{bmatrix} \hat{Y} \end{bmatrix}$

Similarity matrix

threshold: $\geq \sigma$

enumerate

$f_\varphi$ 🔥

feature enhancer

$\mathbf{z}_{v,t}$

$\mathbf{z}_{v,t}$

$\mathbf{z}_{v,t}$

"blueberry ."

$f_{\theta_{TT}}$ ❄️

text encoder

Text tokens

Fused features

$\mathbf{z}_{v,t}$

Count $\hat{y}$

: Added to GroundingDINO

Loss (same as GroundingDINO's but with center points $c_i$ instead of boxes)

$$\lambda_{loc} \sum_{i=1}^{l} |\hat{c}_i - c_i| + \lambda_{cls} \text{FocalLoss}(\hat{\mathbf{Y}}, T)$$

Exemplar Extraction Module

# Training – Dataset

- Trained on open-world object counting dataset FSC-147 [1] with text and visual exemplars.

- Text encoder and image encoder frozen during finetuning.



2.jpg   3.jpg   4.jpg   5.jpg   6.jpg   7.jpg   •••   7705.jpg

21.jpg   22.jpg   23.jpg   24.jpg   26.jpg   27.jpg   7714.jpg

1.  Viresh Ranjan, Udbhav Sharma, Thu Nguyen, and Minh Hoai. Learning to count everything. In *Proc. CVPR*, 2021.

# Results – Qualitative



*From FSC-147 test set*

**Dataset 1**

# Results – Qualitative Continued



"car"

"the world's greatest magicians"

"the beautiful butterfly wall stickers"

Pred: 73, GT: 73

Pred: 7, GT: 7

Pred: 8, GT: 8

*From CARPK*

*From CountBench*

**Dataset 2**

**Dataset 3**

***Zero-shot results with no fine-tuning***

# Results – Qualitative Continued



*From real-world application of trying to understand the influence of climate change on seabird populations. Zero-shot, no fine-tuning.*

# Results – Quantitative

- CountGD achieves SOTA for open-world object counting. *lower is better*.

# App Demo

[CountGD_Multi-Modal_Open-World_Counting - a Hugging Face Space by nikigoli](#)

[https://www.robots.ox.ac.uk/~vgg/research/countgd/](#)

# Thank you!



https://www.robots.ox.ac.uk/~vgg/research/countgd/