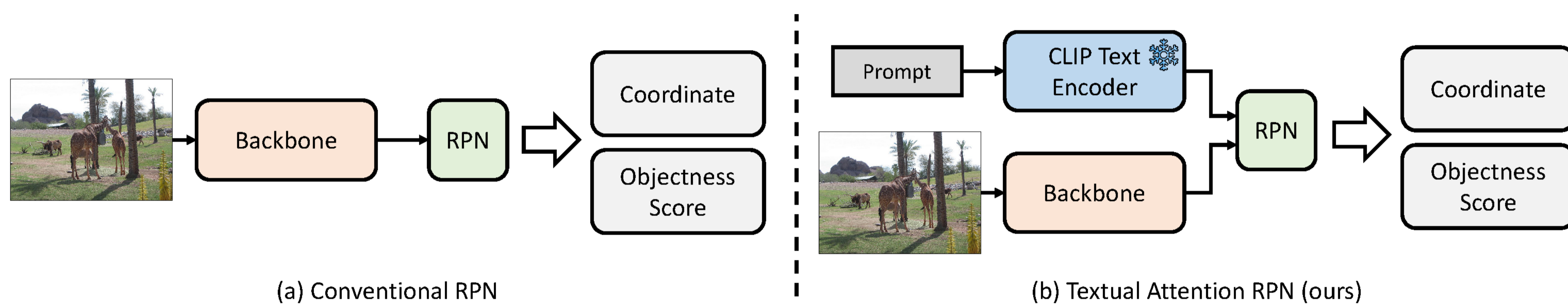


Textual Attention RPN for Open-Vocabulary Object Detection

TL; DR: Make powerful Region Proposal Network (RPN) with textual and visual feature fusion for open-vocabulary object detection

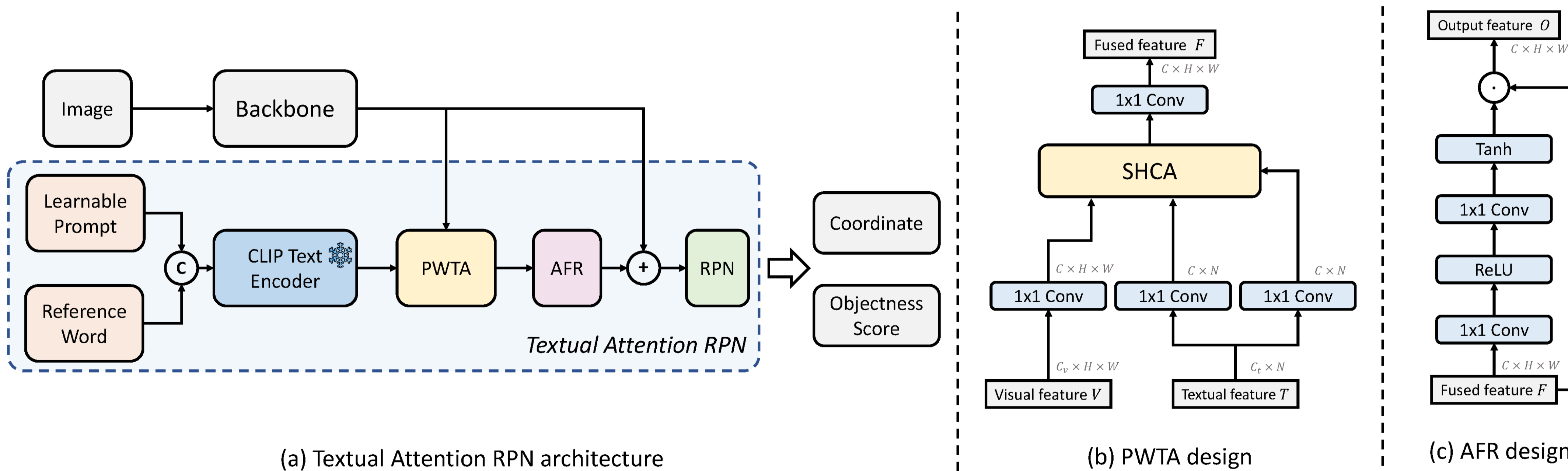
Introduction

- **Open-Vocabulary Object Detection (OVD)** is a method that allows a model to identify objects it hasn't specifically seen or trained on before, using **external textual information**.
- However, **conventional RPNs focus solely on visual features**, which limits their effectiveness for novel objects in OVD..
- Our proposed **Textual Attention RPN (TA-RPN)** integrates both visual and textual features from CLIP's text encoder, employing prompt learning to enhance localization.
- TA-RPN **improves proposal generation performance** through pixel-wise attention and prompt learning.



Method

- TA-RPN integrates information from CLIP's text encoder using **Pixel-Wise Textual Attention (PWTA)**, **Adaptive Feature Refinement (AFR)**, and **prompt learning**.
- **Pixel-Wise Textual Attention (PWTA):** Fuses visual and textual features at each pixel to capture detailed, context-aware information across the entire image.
- **Adaptive Feature Refinement (AFR):** Uses an **attention map** to enhance visual relevance **while preventing textual information from overwhelming the localization process**.
- **Prompt Learning:** Employs dynamic prompts to tailor textual features for effective localization of novel objects.



Experiments

- TA-RPN is built on the **BARON [1] framework** with Faster R-CNN [2] and ResNet50-C4/FPN initialized with SoCo pre-trained weights.
- CLIP (ViT-B/32) is used to generate text embeddings, which are fused with visual features to enhance proposal generation.
- On **COCO**, TA-RPN achieves a **36.1 AP** for novel classes, surpassing existing methods like BARON by **2.1 AP**.
- On **LVIS**, TA-RPN improves the **average precision for rare classes (AP_r) to 21.5**.
- Each module (PWTA, AFR, and prompt learning) was evaluated independently, with all components together resulting in a **4.0 AP increase for novel classes**.
- **The highest average recall (ARN@100) for novel classes** was achieved when all modules were active.

#	PWTA	AFR	Prompt learning	AP ₅₀ ^N	AP ₅₀ ^B	ARN@100
1	X	X	X	32.1	50.7	35.5
2	✓	X	X	34.3 (+2.2)	53.5	35.9 (+0.4)
3	✓	✓	X	35.2 (+3.1)	53.1	36.1 (+0.6)
4	✓	X	✓	34.7 (+2.6)	48.8	35.3 (-0.2)
5	✓	✓	✓	36.1 (+4.0)	53.1	36.4 (+0.9)

Ablation study

Method	Backbone	Extra Data	AP ₅₀ ^N	AP ₅₀ ^B	AP ₅₀
OVR-CNN	ResNet50-C4	✓	22.8	46.0	39.9
ViLD	ResNet50-FPN	X	27.6	59.5	51.3
RegionCLIP	ResNet50-C4	✓	26.8	54.8	50.4
Detic	ResNet50-C4	✓	27.8	47.1	45.0
VLDet	ResNet50-C4	✓	32.0	50.6	45.8
F-VLM	ResNet50-FPN	X	28.0	-	39.6
OADP	ResNet50-FPN	X	30.0	53.3	47.2
CoDet	ResNet50-C4	X	30.6	52.3	46.6
ProxyDet	ResNet50-C4	X	30.4	52.6	46.8
BARON	ResNet50-FPN	X	34.0	60.4	53.5
Ours	ResNet50-C4	X	36.1	53.1	43.3

Performance comparison on the COCO dataset

Method	Backbone	Object Detection				Instance Segmentation			
		AP _r	AP _c	AP _f	AP	AP _r	AP _c	AP _f	AP
ViLD	ResNet50	16.1	20.0	28.3	22.5	16.3	21.2	31.6	24.4
RegionCLIP	ResNet50	17.1	27.4	34.0	28.2	-	-	-	-
Detic	ResNet50	-	-	-	-	17.8	26.3	31.6	26.8
DetPro	ResNet50	20.8	27.8	32.4	28.4	19.8	25.6	28.9	25.9
F-VLM	ResNet50	-	-	-	-	18.6	24.0	26.9	24.2
ProxyDet	ResNet50	-	-	-	-	18.9	-	-	30.1
BARON*	ResNet50	20.4	30.9	33.4	30.1	19.8	28.6	30.2	27.7
Ours	ResNet50	21.5	30.4	33.7	30.2	20.7	28.6	30.3	27.9

Performance comparison on the LVIS dataset

[1] Wu, Size, et al. "Aligning bag of regions for open-vocabulary object detection." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2023.

[2] Ren, Shaoqing, et al. "Faster R-CNN: Towards real-time object detection with region proposal networks." *IEEE transactions on pattern analysis and machine intelligence* 39.6 (2016): 1137-1149.