



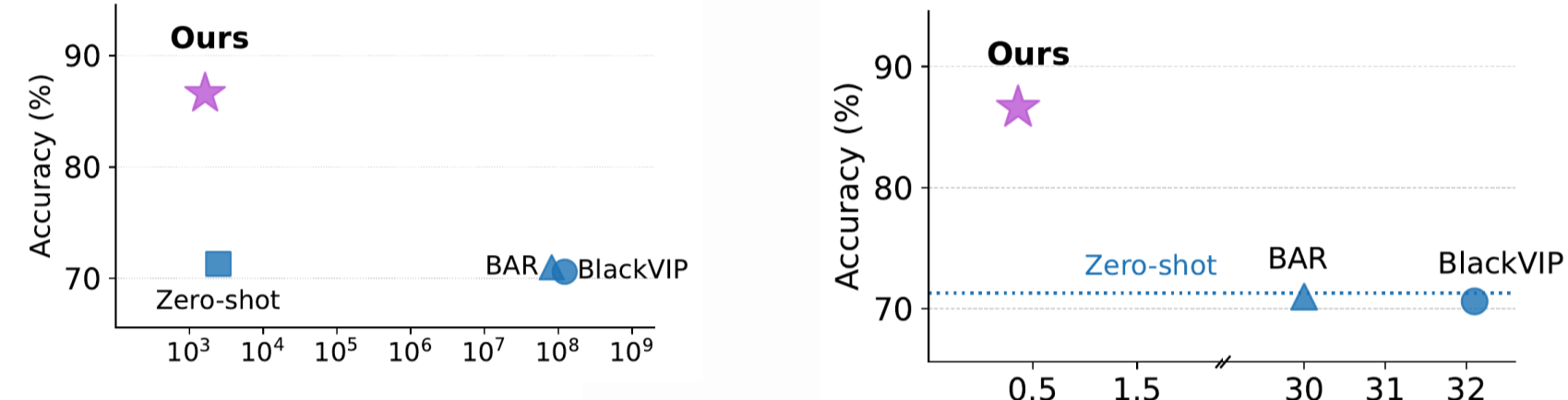
# Prime Once, then Reprogram Locally: An Efficient Alternative to Black-Box Service Model Adaptation

Yunbei Zhang<sup>1</sup>, Chengyi Cai<sup>2</sup>, Feng Liu<sup>2</sup>, Jihun Hamm<sup>1</sup>  
<sup>1</sup>Tulane University <sup>2</sup>University of Melbourne

## Motivation

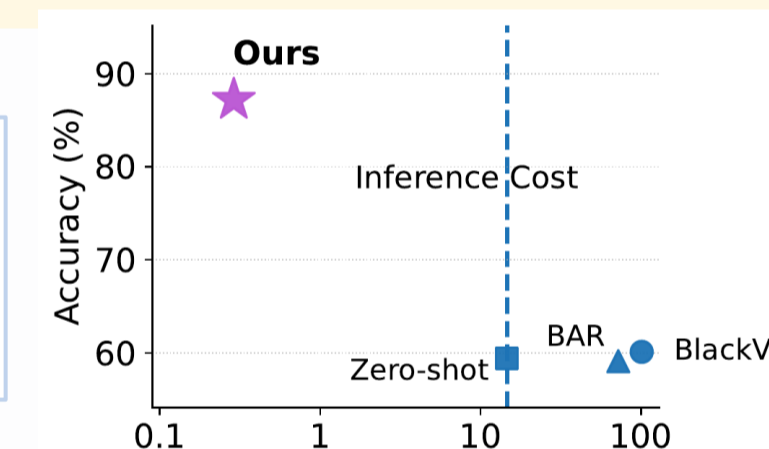
**Black-box reprogramming spends the budget in the wrong place.**

- BAR and BlackVIP optimize visual prompts through Zeroth-Order Optimization (ZOO).
- ZOO needs repeated API calls per image, noisy gradient estimates, and service access again at inference time.
- Modern robust APIs can be insensitive to small input perturbations, so costly search may not move the prediction.



**Question: can we use a service API once, then adapt and deploy locally?**

$$\nabla \ell(\mathbf{P}) \approx \frac{d}{q} \sum_{i=1}^q \frac{\ell(\mathbf{P} + \mu u_i) - \ell(\mathbf{P})}{\mu} u_i$$



Each loss evaluation routes back through the service model.

## Design Goal

- 1x** service pass for priming
- 0** service calls at inference
- Local** gradients after priming

<b>Black Box</b>	input / prediction only
<b>Grey Box</b>	features or embeddings
<b>White Box</b>	gradients or parameters

Access	Method	MR	Input Access	Output Access	VM	VLM	Prompt	Encoder Flexibility	Cost-free Inference
	LFA [42]	x	Input	Features	x	x	-	-	x
	CBBT [13]	x	Embedding	Features	x	x	L	-	x
	LaSP [20]	x	Input	Features	x	x	V	-	x
	BPT-VLM [73]	x	Input	Predictions	x	x	L&V	VIT Only	x
	CrPT [54]	x	Embedding	Predictions	x	x	L	-	x
	BAR [7]	x	Input	Predictions	x	x	V	-	x
	BlackVIP [41]	x	Input	Predictions	x	x	V	-	x
	LLM-Opt [13]	x	Input	Predictions	x	x	L	-	x
	AReS (Ours)	x	Input	Predictions	x	x	V	-	x

- Stay in strict Black Box access for the service API.
- Prime once, then learn and deploy the visual prompt locally.
- Use the same primitive for VMs and VLMs.

## AReS: Prime Once, Then Reprogram Locally

### 1. Query Service

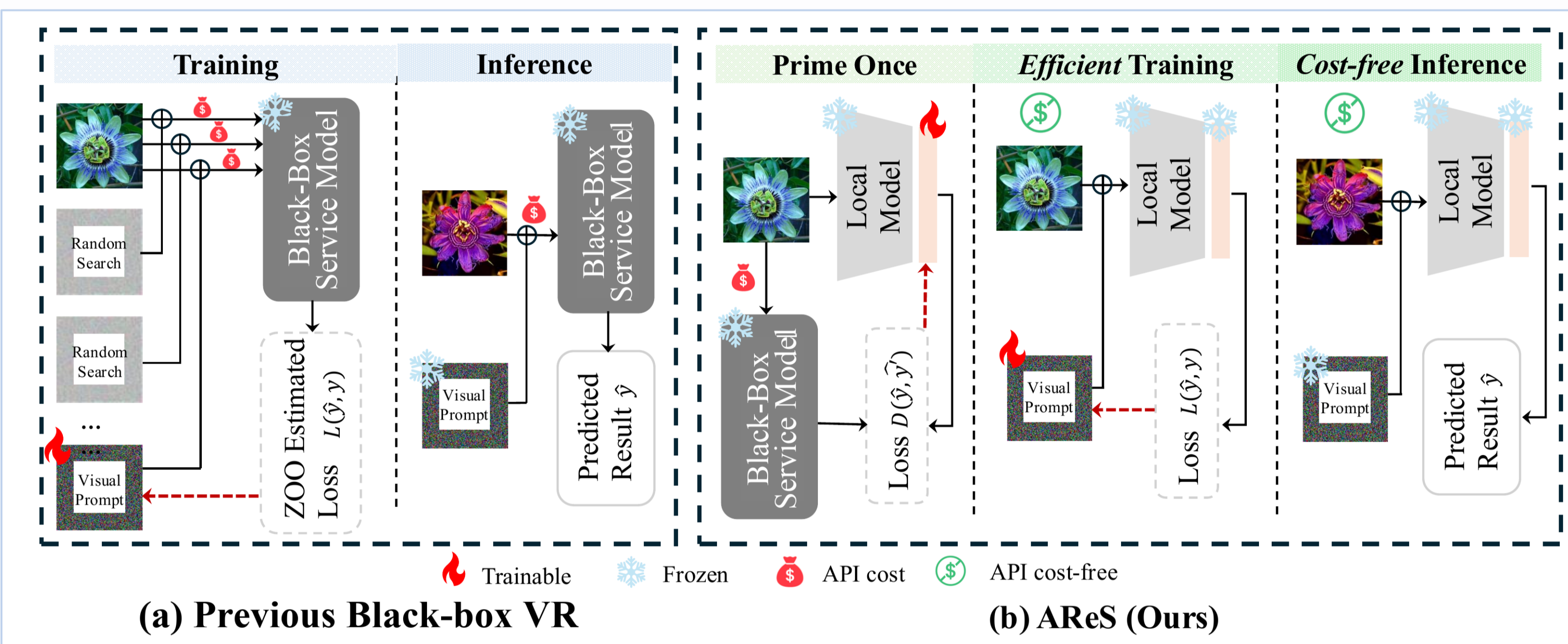
Collect service probabilities once per target-domain training image.

### 2. Prime Local

Freeze the local encoder and train a lightweight linear layer.

### 3. Local VR

Learn the visual prompt with exact White Box gradients on the primed local model.



AReS moves adaptation and inference to the local proxy after one-time priming.

### Priming loss

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}_p(\mathcal{F}_L(x_i; \theta), p_S(x_i))$$

### Local VR loss

$$\mathbf{P}^* = \arg \min_{\mathbf{P}} \mathbb{E}_{(x,y) \sim \mathbb{D}^T} [\ell(g_{out}(\mathcal{F}_L(g_{in}(x; \mathbf{P}); \theta^*)), y)]$$

Method details are the two optimization objectives: one service pass for priming, then all prompt learning with local gradients.

## Main Results: Accuracy and Efficiency

AReS improves accuracy while replacing repeated service-model optimization with one-time priming.

Table 2. VLM / VM accuracy, API calls, and training time

Method	Flowers	DTD	UCF	Food	GTSRB	EuroSAT	Pets	Cars	SUN	SVHN	Avg.	#API (M)	Time (h)
VR (glass-box)	86.8	62.0	74.0	81.6	65.1	90.9	90.7	66.2	66.7	60.2	74.4	16.82	9.8
Zero-shot	71.3	43.9	66.9	<b>85.9</b>	21.0	47.9	<b>89.1</b>	65.2	62.6	17.9	57.2	0.12	0
BAR	71.0	46.8	64.2	84.4	21.5	77.3	88.4	63.0	62.4	34.6	61.4	612.84	185.6
BlackVIP	70.6	45.3	<b>68.7</b>	<b>85.9</b>	21.3	73.3	<b>89.1</b>	<b>65.4</b>	<b>64.5</b>	44.4	62.9	754.20	197.5
LLM-Opt	67.6	45.0	59.9	78.5	21.2	48.0	87.7	56.2	60.3	20.2	54.5	5011.32	5.1
<b>AReS</b>	<b>86.6</b>	<b>48.2</b>	<b>67.1</b>	<b>85.9</b>	<b>39.4</b>	<b>85.7</b>	<b>88.9</b>	<b>43.2</b>	<b>62.8</b>	<b>63.2</b>	<b>65.4</b>	<b>0.02</b>	<b>3.7</b>
<b>AReS-MS</b>	<b>86.6</b>	<b>48.2</b>	<b>67.1</b>	<b>85.9</b>	<b>39.4</b>	<b>85.7</b>	<b>88.9</b>	<b>65.2</b>	<b>62.8</b>	<b>63.2</b>	<b>69.3</b>	<b>0.06</b>	<b>3.7</b>

Table 4. Modern API case study

Method	LLaVA			GPT-4o			Clarifai		
	Acc. ↑	Train	Infer. Total (#) ↓	Acc. ↑	Train	Infer. Total (\$) ↓	Acc. ↑	Train	Infer. Total (\$) ↓
glass-box	91.3	-	-	-	-	-	-	-	-
Zero-shot	40.1	-	8100	59.4	-	14.6	14.6	-	-
BAR	34.1	~ 10 <sup>6</sup>	8100	59.1	57.6	14.6	72.2	68.1	38.4
BlackVIP	39.4	~ 10 <sup>7</sup>	8100	60.1	86.4	14.6	101.0	72.1	57.6
<b>AReS</b>	<b>73.1</b>	160	0	<b>87.2</b>	0.3	0	<b>0.3</b>	<b>83.2</b>	0.2

Efficiency reading: the accuracy gains come with cost-free inference and no repeated API perturbation loop.

**User**

You are an expert flower classifier for the Oxford Flowers-102 dataset. Analyze this image and identify the flower species from these categories: [pink primrose, ...]. Consider petal shape, color, structure, and distinctive features. Provide your response in exactly this format: First paragraph: Briefly describe the key visible characteristics and explain your classification reasoning in 1-2 sentences. Second paragraph: "Predicted Class: [Flower Name]"

Clean

Predicted Class: **Artichoke**

**User**

You are an expert flower classifier for the Oxford Flowers-102 dataset. Analyze this image and identify the flower species from these categories: [pink primrose, ...]. Consider petal shape, color, structure, and distinctive features. Provide your response in exactly this format: First paragraph: Briefly describe the key visible characteristics and explain your classification reasoning in 1-2 sentences. Second paragraph: "Predicted Class: [Flower Name]"

Prompted

Predicted Class: **Artichoke**

AReS is less dependent on service-model perturbation sensitivity because optimization happens on the primed local model.

## Theoretical Results

- Priming is useful when it makes the local proxy faithful to the service model.
- Assume ( $\epsilon$ ) –faithful priming on target-domain inputs.
- With cross-entropy Lipschitzness, logit closeness transfers to risk.

$$\mathbb{E}_{\mathbb{D}^T} |z_S^*(x, \mathbf{Q}^*) - z_L^*(x, \mathbf{P}^*)|_1 \leq \epsilon$$

$$\mathcal{R}_L(\mathcal{D}^T, \mathbf{P}^*) - \epsilon \leq \mathcal{R}_S(\mathcal{D}^T, \mathbf{Q}^*) \leq \mathcal{R}_L(\mathcal{D}^T, \mathbf{P}^*)$$

- The epsilon-faithful assumption is not free; it is exactly the goal of the priming stage.

- The mechanism is constructive: one service pass trains a lightweight head on a frozen local encoder.

- The downstream label space need not match the service label space, so priming is preparatory rather than ordinary final-task distillation.

**Insight: faithful priming lets a stable White Box local route approximate the reprogrammed service model.**

The theorem explains why API efficiency and stability improve together.

## Conclusion and Discussion

**Take-home: AReS turns black-box service adaptation into one-time transfer plus local, cost-free deployment.**

### Key conclusions

- Modern APIs are fragile under perturbation-based ZOO.

- AReS spends API calls once, then optimizes locally.

- Future work: make priming more faithful under larger domain gaps.