



# FLiP: Towards Comprehensive and Reliable Evaluation of Federated Prompt Learning

Dongping Liao<sup>1</sup>, Xitong Gao<sup>2,3</sup>, Chengzhong Xu<sup>1</sup>

<sup>1</sup>University of Macau, <sup>2</sup>Shenzhen Institutes of Advanced Technology,

<sup>3</sup>Shenzhen University of Advanced Technology



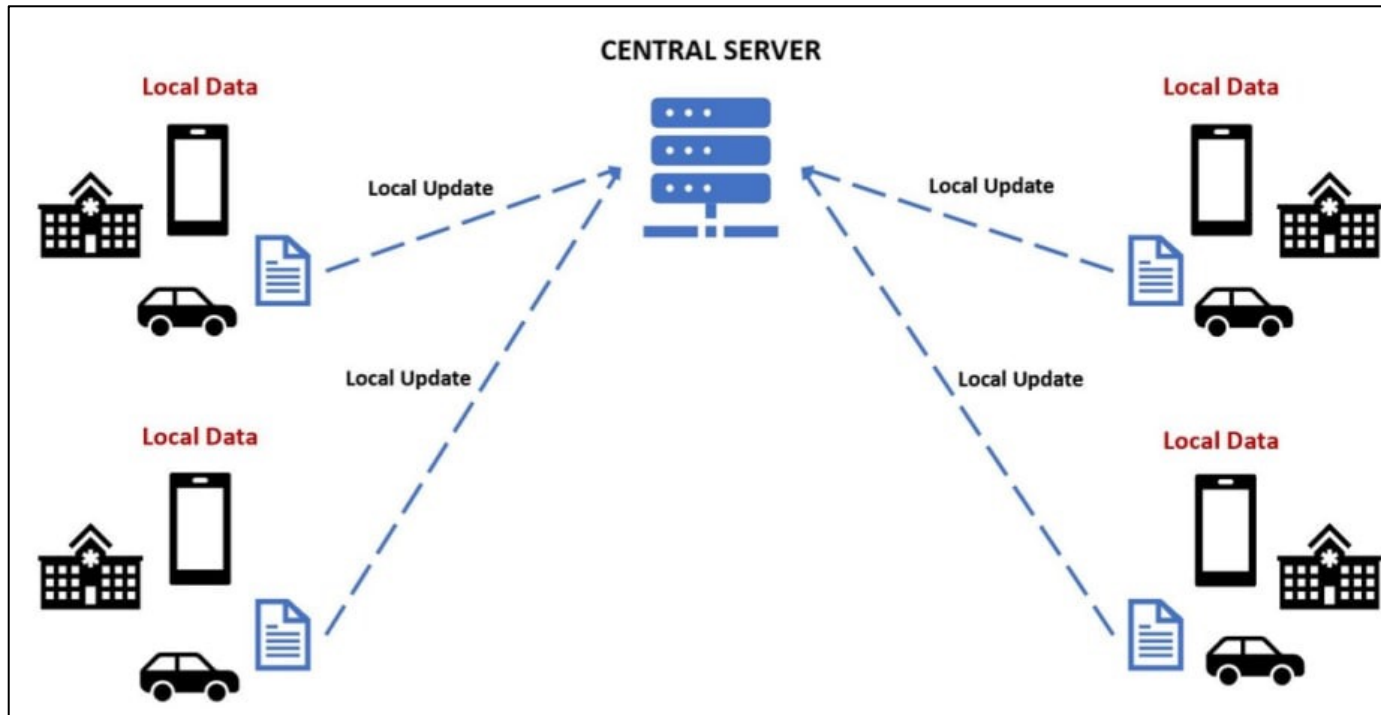
澳門大學  
UNIVERSIDADE DE MACAU  
UNIVERSITY OF MACAU



DEPARTMENT OF  
COMPUTER AND INFORMATION SCIENCE

# Background

## Federated Learning



**Figure:** A conceptual illustration of a federated learning system.

# Background

## Limitations of Current Federated Prompt Learning Evaluation Approaches

- A **systematic understanding** of their underlying mechanisms and **principled guidelines** for deploying these techniques in different FL scenarios remain absent.
- Existing **inconsistent experimental protocols, limited evaluation scenarios**, and the lack of the proper assessment of **centralized PL methods** have obscured the essence of these algorithms.

To close above gaps, we introduce a **comprehensive benchmark**, named FLiP, to achieve standardized FPL evaluation.



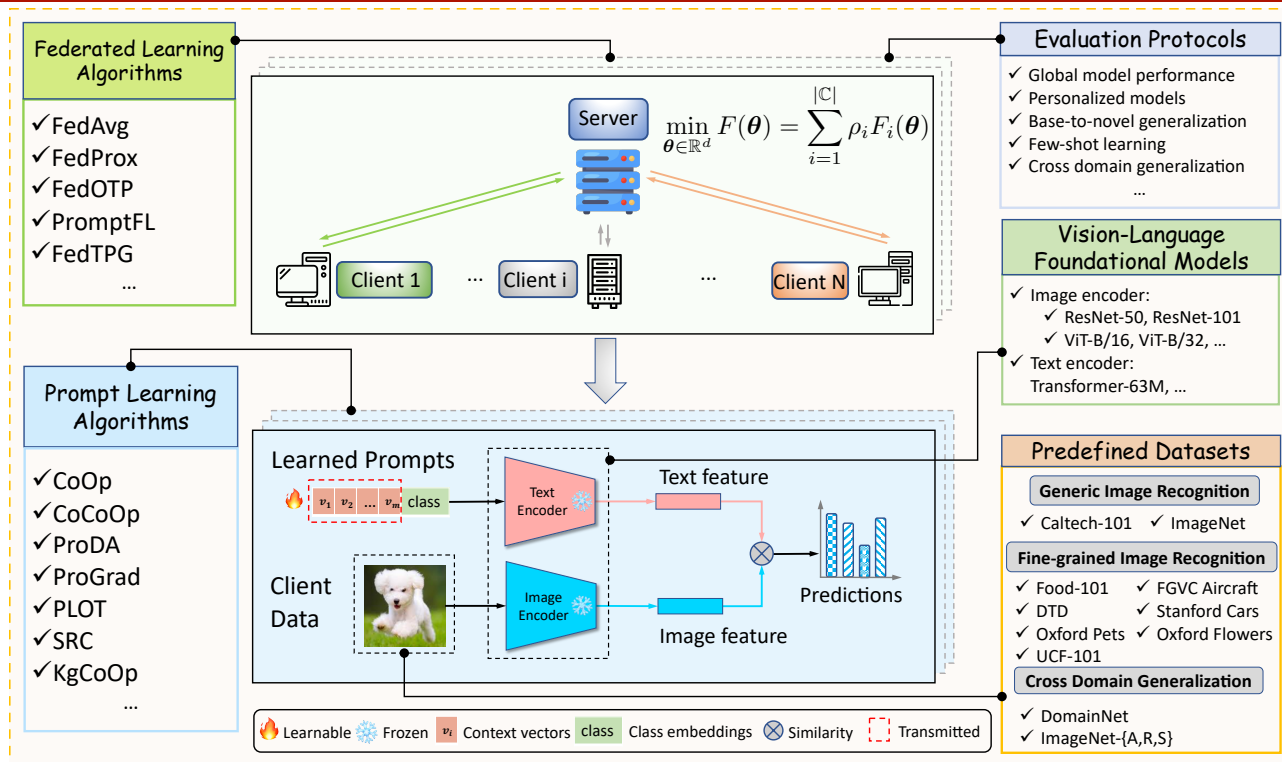
# Research Questions

- **Q1:** How effective are **global and personalized prompt models** obtained by FPL methods?
- **Q2:** How do **various real-world data distribution shifts** impact FPL?
- **Q3:** What are **best practices** for deploying these techniques in different FL scenarios?
- **Q4:** What can we learn about **cost-effectiveness** from different FPL methods?
- **Q5:** What are the **remaining challenges in FPL?**



# Benchmark Design

## The System Design of the FLiP Benchmark



### Core Features

- Decoupled FL and PL modules
- Standardized data interfaces
- Various VL foundational models
- Extensible open-source codebase

### Benchmark Statistics

- 3 evaluation protocols
- 6 FPL scenarios
- 12 open datasets
- 13 PL and FPL methods



# Benchmark Results

## Global Model Performance

Global $\alpha_g$	Caltech	DTD	Aircraft	Food	Cars	Flowers	Pets	UCF	Avg.	#
<b>ZS-CLIP</b>	86.0	41.7	16.6	77.9	55.5	65.3	85.7	61.5	61.3	-
<b>PromptFL</b>	91.5 $\pm$ 0.5	57.6 $\pm$ 1.3	22.8 $\pm$ 0.4	79.2 $\pm$ 0.1	62.0 $\pm$ 0.4	<b>84.0<math>\pm</math>1.7</b>	<b>89.4<math>\pm</math>0.5</b>	70.1 $\pm$ 0.8	<u>69.6</u>	-
<b>FedOTP</b>	91.8 $\pm$ 0.1	58.0 $\pm$ 0.8	21.9 $\pm$ 0.4	78.7 $\pm$ 0.1	62.8 $\pm$ 0.2	83.3 $\pm$ 0.6	89.1 $\pm$ 0.1	69.4 $\pm$ 0.6	69.4	3
<b>FedTPG</b>	90.2 $\pm$ 0.1	56.8 $\pm$ 1.0	19.0 $\pm$ 1.2	79.3 $\pm$ 0.2	60.7 $\pm$ 0.2	78.0 $\pm$ 1.6	89.0 $\pm$ 0.6	68.3 $\pm$ 0.2	67.6	1
<b><i>f</i>-CoCoOp</b>	91.7 $\pm$ 0.3	54.7 $\pm$ 1.0	17.9 $\pm$ 4.5	79.3 $\pm$ 0.1	60.7 $\pm$ 0.5	76.4 $\pm$ 0.7	89.1 $\pm$ 0.1	68.0 $\pm$ 0.9	67.2	2
<b><i>f</i>-PLOT</b>	91.6 $\pm$ 0.3	<b>58.3<math>\pm</math>1.4</b>	21.7 $\pm$ 0.5	78.3 $\pm$ 0.2	60.7 $\pm$ 0.8	83.4 $\pm$ 0.6	88.9 $\pm$ 0.5	69.7 $\pm$ 0.4	69.1	2
<b><i>f</i>-ProDA</b>	91.6 $\pm$ 0.3	57.2 $\pm$ 1.1	<b>23.1<math>\pm</math>0.7</b>	79.1 $\pm$ 0.2	62.3 $\pm$ 0.5	<b>84.0<math>\pm</math>0.8</b>	<u>89.3<math>\pm</math>0.4</u>	70.2 $\pm$ 1.1	<u>69.6</u>	<b>5</b>
<b><i>f</i>-ProGrad</b>	90.7 $\pm$ 0.2	57.1 $\pm$ 1.0	21.7 $\pm$ 0.3	<b>79.5<math>\pm</math>0.1</b>	60.5 $\pm$ 0.7	83.4 $\pm$ 0.4	89.1 $\pm$ 0.2	<u>70.3<math>\pm</math>0.3</u>	69.1	2
<b><i>f</i>-PromptSRC</b>	<b>92.0<math>\pm</math>0.8</b>	57.8 $\pm$ 0.3	21.2 $\pm$ 0.4	78.6 $\pm$ 0.4	62.4 $\pm$ 0.2	83.6 $\pm$ 0.1	89.2 $\pm$ 0.7	<u>70.3<math>\pm</math>1.0</u>	69.4	<u>4</u>
<b><i>f</i>-KgCoOp</b>	91.8 $\pm$ 0.2	<u>58.2<math>\pm</math>0.8</u>	<u>23.0<math>\pm</math>0.1</u>	<u>79.4<math>\pm</math>0.2</u>	61.7 $\pm$ 0.7	<u>83.9<math>\pm</math>0.5</u>	<b>89.4<math>\pm</math>0.2</b>	<b>70.4<math>\pm</math>0.7</b>	<b>69.7</b>	<b>5</b>

PromptFL, combining CoOp and FedAvg, serves as a simple yet effective baseline and is competitive on fine-grained datasets like Oxford-Pets and Flowers.



# Benchmark Results

## Personalized Model Performance

Personal $\alpha_p$	Caltech	DTD	Aircraft	Food	Cars	Flowers	Pets	UCF	Avg.	#
<b>ZS-CLIP</b>	86.0	41.7	16.6	77.9	55.5	65.3	85.7	61.5	61.3	-
<b>PromptFL</b>	91.5 $\pm$ 0.4	69.5 $\pm$ 4.1	33.8 $\pm$ 0.1	82.1 $\pm$ 0.4	67.7 $\pm$ 0.6	89.7 $\pm$ 0.2	89.9 $\pm$ 0.6	77.5 $\pm$ 1.5	75.2	0
<b>FedOTP</b>	91.9 $\pm$ 0.4	73.8 $\pm$ 1.4	36.1 $\pm$ 0.5	82.0 $\pm$ 0.7	68.1 $\pm$ 1.7	89.6 $\pm$ 0.3	89.5 $\pm$ 1.2	80.7 $\pm$ 1.0	76.5	5
<b>FedTPG</b>	89.9 $\pm$ 0.4	66.8 $\pm$ 0.6	31.6 $\pm$ 0.3	81.9 $\pm$ 0.3	66.2 $\pm$ 0.2	86.0 $\pm$ 0.4	88.9 $\pm$ 0.8	78.0 $\pm$ 0.7	73.7	1
<b>FedPGP</b>	91.8 $\pm$ 0.4	68.0 $\pm$ 0.5	35.2 $\pm$ 0.4	81.0 $\pm$ 0.2	66.7 $\pm$ 0.6	84.1 $\pm$ 1.6	87.4 $\pm$ 0.3	77.8 $\pm$ 0.5	74.0	3
<b>PromptFolio</b>	91.6 $\pm$ 0.3	70.2 $\pm$ 0.4	34.6 $\pm$ 0.5	82.0 $\pm$ 0.3	67.4 $\pm$ 0.2	89.2 $\pm$ 0.5	89.4 $\pm$ 0.2	79.2 $\pm$ 0.6	75.5	4
<b>DP-FPL</b>	90.2 $\pm$ 0.4	65.2 $\pm$ 0.5	28.5 $\pm$ 0.5	80.1 $\pm$ 0.6	64.5 $\pm$ 1.4	78.6 $\pm$ 0.6	82.4 $\pm$ 0.8	72.2 $\pm$ 0.7	70.2	0
<i>f</i> -CoCoOp	91.8 $\pm$ 0.4	70.3 $\pm$ 3.0	34.0 $\pm$ 1.7	81.8 $\pm$ 0.6	67.4 $\pm$ 0.7	86.4 $\pm$ 1.9	89.5 $\pm$ 0.9	77.1 $\pm$ 0.9	74.8	3
<i>f</i> -PLOT	91.7 $\pm$ 0.4	71.3 $\pm$ 3.1	34.0 $\pm$ 0.8	81.4 $\pm$ 1.2	67.9 $\pm$ 1.1	89.3 $\pm$ 0.5	88.6 $\pm$ 0.3	79.6 $\pm$ 1.8	75.5	5
<i>f</i> -ProDA	91.7 $\pm$ 0.9	69.7 $\pm$ 2.6	34.7 $\pm$ 0.6	82.1 $\pm$ 1.3	67.8 $\pm$ 1.3	89.3 $\pm$ 0.6	89.9 $\pm$ 0.5	77.9 $\pm$ 1.5	75.4	6
<i>f</i> -ProGrad	91.7 $\pm$ 0.6	69.2 $\pm$ 0.3	32.9 $\pm$ 0.7	81.6 $\pm$ 0.8	67.0 $\pm$ 1.0	88.8 $\pm$ 0.9	89.5 $\pm$ 0.7	77.0 $\pm$ 1.5	74.7	1
<i>f</i> -PromptSRC	91.7 $\pm$ 0.4	69.3 $\pm$ 1.3	32.4 $\pm$ 2.3	81.8 $\pm$ 1.1	67.6 $\pm$ 1.1	89.2 $\pm$ 1.7	89.3 $\pm$ 1.7	78.2 $\pm$ 1.1	74.9	2
<i>f</i> -KgCoOp	91.6 $\pm$ 0.3	68.6 $\pm$ 2.7	31.3 $\pm$ 0.5	81.4 $\pm$ 0.7	67.1 $\pm$ 1.4	88.9 $\pm$ 0.9	89.9 $\pm$ 0.3	76.9 $\pm$ 0.9	74.5	1

FedOTP generally outperforms other personalized methods, highlighting the efficacy of distribution alignment in adapting to personalized data.





# Benchmark Results

## Base-to-Novel Class Generalization

	Caltech			Aircraft			Cars			Flowers			Avg.			
Metric	$\alpha_b$	$\alpha_n$	$\alpha_h$	$\alpha_b$	$\alpha_n$	$\alpha_h$	$\alpha_b$	$\alpha_n$	$\alpha_h$	$\alpha_b$	$\alpha_n$	$\alpha_h$	$\alpha_b$	$\alpha_n$	$\alpha_h$	#
<b>ZS-CLIP</b>	88.2	92.6	90.3	19.6	24.7	21.8	59.5	68.1	63.5	77.2	71.0	73.9	61.1	64.1	62.4	-
<b>PromptFL</b>	92.8 $\pm$ 0.8	92.7 $\pm$ 0.7	92.6 $\pm$ 0.2	20.9 $\pm$ 0.4	24.7 $\pm$ 0.5	22.6 $\pm$ 0.2	63.0 $\pm$ 0.7	67.6 $\pm$ 0.7	65.2 $\pm$ 0.1	79.9 $\pm$ 2.9	69.3 $\pm$ 0.7	74.2 $\pm$ 1.0	64.1	63.8	63.8	-
<b>FedOTP</b>	93.1 $\pm$ 0.2	93.7 $\pm$ 0.4	93.4 $\pm$ 0.1	21.0 $\pm$ 0.8	23.7 $\pm$ 1.0	22.2 $\pm$ 0.8	62.1 $\pm$ 0.1	66.0 $\pm$ 0.7	64.0 $\pm$ 0.4	81.0 $\pm$ 0.6	68.7 $\pm$ 1.9	74.3 $\pm$ 1.2	64.3	63.0	63.5	2
<b>FedTPG</b>	93.6 $\pm$ 0.4	90.0 $\pm$ 0.6	91.8 $\pm$ 0.4	19.5 $\pm$ 0.3	22.7 $\pm$ 0.5	21.0 $\pm$ 0.5	66.4 $\pm$ 0.2	67.7 $\pm$ 0.2	<b>67.0<math>\pm</math>0.2</b>	75.0 $\pm$ 0.3	66.0 $\pm$ 0.5	70.2 $\pm$ 0.4	63.6	61.6	62.5	2
<b>FedPGP</b>	93.2 $\pm$ 0.4	92.7 $\pm$ 1.1	92.9 $\pm$ 0.5	19.8 $\pm$ 1.0	19.2 $\pm$ 2.1	19.5 $\pm$ 0.8	63.7 $\pm$ 1.4	67.8 $\pm$ 0.7	65.7 $\pm$ 0.9	80.7 $\pm$ 1.4	66.8 $\pm$ 1.7	73.1 $\pm$ 1.2	64.3	61.6	62.8	2
<b>f-CoCoOp</b>	92.7 $\pm$ 0.8	93.6 $\pm$ 0.6	93.1 $\pm$ 0.2	18.0 $\pm$ 2.1	17.1 $\pm$ 2.4	17.2 $\pm$ 2.2	62.8 $\pm$ 0.6	66.7 $\pm$ 0.3	64.7 $\pm$ 0.2	79.4 $\pm$ 1.4	70.7 $\pm$ 1.8	74.8 $\pm$ 0.6	63.2	62.0	62.4	2
<b>f-PLOT</b>	93.4 $\pm$ 0.6	93.5 $\pm$ 1.0	<b>93.5<math>\pm</math>0.8</b>	19.0 $\pm$ 0.8	23.4 $\pm$ 0.4	21.0 $\pm$ 0.4	62.4 $\pm$ 0.5	65.4 $\pm$ 1.4	63.8 $\pm$ 0.7	78.7 $\pm$ 3.0	68.3 $\pm$ 1.3	73.1 $\pm$ 0.6	63.4	62.6	62.8	1
<b>f-ProDA</b>	93.0 $\pm$ 0.4	92.7 $\pm$ 1.1	92.8 $\pm$ 0.4	22.0 $\pm$ 0.5	25.1 $\pm$ 1.1	23.4 $\pm$ 0.6	63.3 $\pm$ 0.7	67.7 $\pm$ 0.4	65.4 $\pm$ 0.2	77.9 $\pm$ 0.7	69.6 $\pm$ 0.5	73.5 $\pm$ 0.1	64.0	63.8	63.8	3
<b>f-ProGrad</b>	93.2 $\pm$ 0.4	93.0 $\pm$ 0.4	93.1 $\pm$ 0.4	21.6 $\pm$ 0.5	25.2 $\pm$ 1.6	23.3 $\pm$ 0.4	64.2 $\pm$ 0.6	67.9 $\pm$ 0.5	66.0 $\pm$ 0.3	80.3 $\pm$ 1.1	70.6 $\pm$ 0.5	<b>75.2<math>\pm</math>0.2</b>	64.7	64.2	<b>64.3</b>	4
<b>f-SRC</b>	90.2 $\pm$ 0.2	93.0 $\pm$ 0.1	91.6 $\pm$ 0.1	21.7 $\pm$ 1.0	24.9 $\pm$ 1.0	23.1 $\pm$ 0.8	61.2 $\pm$ 0.2	67.1 $\pm$ 0.3	64.0 $\pm$ 0.1	78.9 $\pm$ 1.1	70.2 $\pm$ 0.8	74.3 $\pm$ 0.7	62.3	63.8	62.8	3
<b>f-KgCoOp</b>	93.5 $\pm$ 0.4	93.4 $\pm$ 1.0	93.4 $\pm$ 0.5	22.3 $\pm$ 0.7	25.1 $\pm$ 0.5	<b>23.6<math>\pm</math>0.5</b>	63.7 $\pm$ 0.3	67.5 $\pm$ 0.3	65.6 $\pm$ 0.1	79.6 $\pm$ 2.7	68.9 $\pm$ 1.2	73.9 $\pm$ 0.5	64.8	63.7	64.1	3

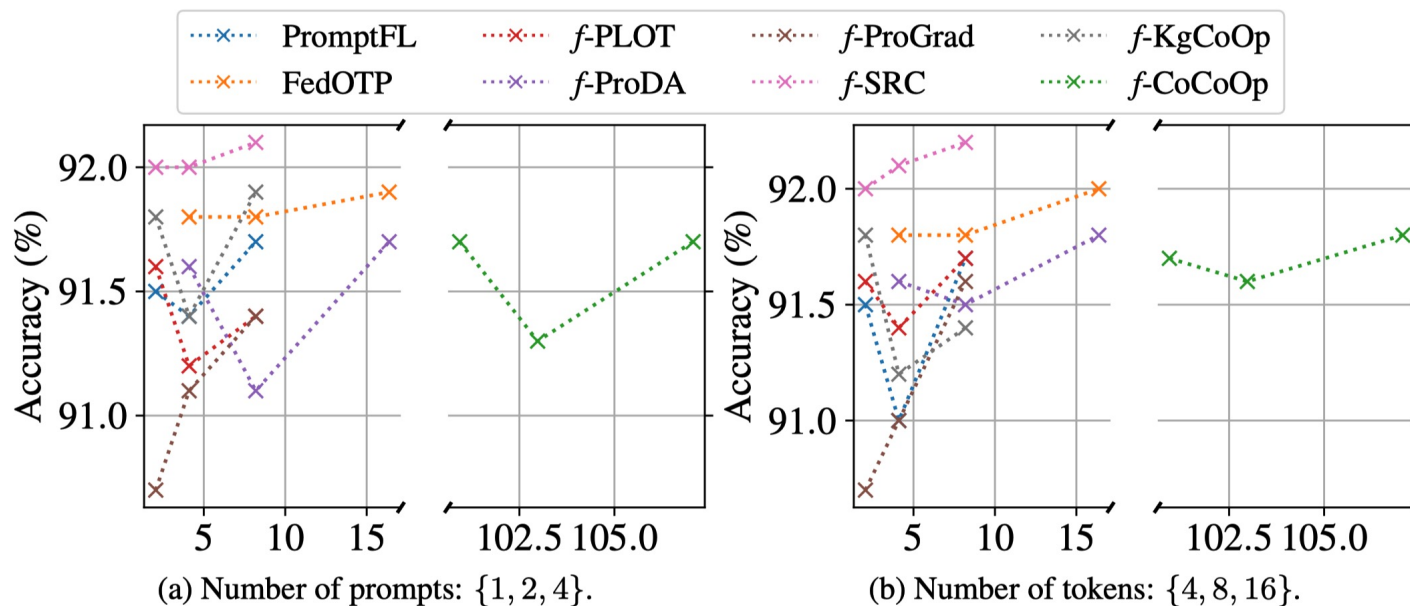
Regularization prevents overfitting and balances base and novel class metrics, without it (PromptFL) is less effective.





# Benchmark Results

## Cost-Performance Tradeoff

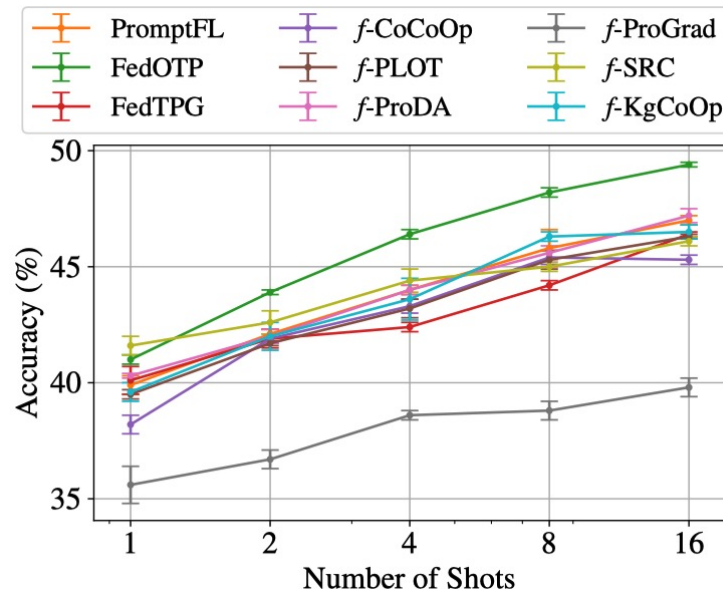


Extra learnable parameters (CoCoOp) do not necessarily improve performance but could be detrimental to computational and communication efficiency.

# Benchmark Results

## Few-shot Accuracy

$k =$	1	2	4	8	16	Avg.	#
<b>PromptFL</b>	39.9 $\pm$ 0.4	42.1 $\pm$ 0.2	44.0 $\pm$ 0.2	45.8 $\pm$ 0.8	47.0 $\pm$ 0.2	43.7	-
<b>FedOTP</b>	41.0 $\pm$ 0.2	<b>43.9<math>\pm</math>0.1</b>	<b>46.4<math>\pm</math>0.2</b>	<b>48.2<math>\pm</math>0.2</b>	<b>49.4<math>\pm</math>0.1</b>	<b>45.8</b>	<b>5</b>
<b>FedTPG</b>	40.1 $\pm$ 0.6	41.9 $\pm$ 0.4	42.4 $\pm$ 0.2	44.2 $\pm$ 0.4	46.4 $\pm$ 0.2	41.8	1
<i>f</i> -CoCoOp	38.2 $\pm$ 0.4	41.9 $\pm$ 0.2	43.3 $\pm$ 0.3	45.4 $\pm$ 0.3	45.3 $\pm$ 0.2	42.8	0
<i>f</i> -PLOT	39.5 $\pm$ 0.2	41.7 $\pm$ 0.1	43.2 $\pm$ 0.4	45.3 $\pm$ 0.4	46.3 $\pm$ 0.1	43.2	0
<i>f</i> -ProDA	40.3 $\pm$ 0.1	42.0 $\pm$ 0.3	44.0 $\pm$ 0.2	45.6 $\pm$ 0.3	47.2 $\pm$ 0.3	43.8	2
<i>f</i> -ProGrad	35.6 $\pm$ 0.8	36.7 $\pm$ 0.4	38.6 $\pm$ 0.2	38.8 $\pm$ 0.4	39.8 $\pm$ 0.4	37.9	0
<b><i>f</i>-SRC</b>	<b>41.6<math>\pm</math>0.4</b>	42.6 $\pm$ 0.5	44.4 $\pm$ 0.5	45.0 $\pm$ 0.2	46.1 $\pm$ 0.2	43.9	3
<i>f</i> -KgCoOp	39.6 $\pm$ 0.4	42.0 $\pm$ 0.6	43.6 $\pm$ 0.9	46.3 $\pm$ 0.2	46.5 $\pm$ 0.3	43.6	1



Methods that use multiple prompts (e.g., FedOTP, f-ProDA, f-SRC) perform best in few-shot scenarios, indicating ensemble helps reduce sample selection bias.

# Benchmark Results

## Feature Shift Heterogeneity

Feature Shift	Clipart	Infograph	Painting	Quickdraw	Real	Sketch	Avg.	#
<b>ZS-CLIP</b>	54.8	40.9	48.8	6.0	77.7	49.3	44.6	-
<b>PromptFL</b>	59.6 $\pm$ 0.2	45.6 $\pm$ 0.3	<u>53.7</u> $\pm$ 0.3	8.9 $\pm$ 0.1	79.8 $\pm$ 0.1	<b>54.2</b> $\pm$ 0.2	<u>48.1</u> $\pm$ 0.1	-
<b>FedOTP</b>	58.4 $\pm$ 0.2	45.2 $\pm$ 0.1	53.4 $\pm$ 0.2	<u>9.0</u> $\pm$ 0.1	79.2 $\pm$ 0.1	53.2 $\pm$ 0.1	47.7 $\pm$ 0.1	1
<b>FedTPG</b>	59.8 $\pm$ 0.1	45.8 $\pm$ 0.3	53.6 $\pm$ 0.2	8.5 $\pm$ 0.2	79.9 $\pm$ 0.3	54.0 $\pm$ 0.3	47.9 $\pm$ 0.3	3
<b><i>f</i>-CoCoOp</b>	<b>60.0</b> $\pm$ 0.1	<b>46.1</b> $\pm$ 0.2	53.0 $\pm$ 0.3	<b>9.1</b> $\pm$ 0.2	79.8 $\pm$ 0.2	<b>54.2</b> $\pm$ 0.2	<u>48.1</u> $\pm$ 0.2	<b>5</b>
<b><i>f</i>-PLOT</b>	58.5 $\pm$ 0.3	44.8 $\pm$ 0.2	53.0 $\pm$ 0.1	<u>9.0</u> $\pm$ 0.4	79.2 $\pm$ 0.1	53.3 $\pm$ 0.1	47.6 $\pm$ 0.1	1
<b><i>f</i>-ProDA</b>	59.5 $\pm$ 0.2	45.6 $\pm$ 0.1	<b>53.8</b> $\pm$ 0.2	<u>9.0</u> $\pm$ 0.2	79.6 $\pm$ 0.1	54.0 $\pm$ 0.3	48.0 $\pm$ 0.1	3
<b><i>f</i>-ProGrad</b>	58.8 $\pm$ 0.2	44.5 $\pm$ 0.2	52.5 $\pm$ 0.1	7.5 $\pm$ 0.2	<u>80.0</u> $\pm$ 0.1	53.0 $\pm$ 0.1	47.3 $\pm$ 0.1	1
<b><i>f</i>-SRC</b>	59.0 $\pm$ 0.1	44.6 $\pm$ 0.4	52.6 $\pm$ 0.1	7.8 $\pm$ 0.1	79.7 $\pm$ 0.1	52.9 $\pm$ 0.1	47.3 $\pm$ 0.1	0
<b><i>f</i>-KgCoOp</b>	<u>59.9</u> $\pm$ 0.1	<u>45.9</u> $\pm$ 0.2	53.6 $\pm$ 0.1	8.8 $\pm$ 0.1	<b>80.2</b> $\pm$ 0.1	<u>54.1</u> $\pm$ 0.1	<b>48.2</b> $\pm$ 0.1	<u>4</u>

It remains challenging to counteract the adverse impact of feature shift for all evaluated FPL methods.

# Benchmark Results

## Cross-domain Generalization

Table 8: Comparing the cross-domain performance of FPL methods. Here, the source domain is ImageNet (IN), and the target domains are ImageNet-A(dversarial), -R(endition), and -S(ketch), respectively denoted as IN-A, IN-R, IN-S.

Cross-domain $\alpha_{\text{IN} \rightarrow \cdot}$	IN-A	IN-R	IN-S	Avg.	#
<b>ZS-CLIP</b>	21.7	56.1	33.4	37.1	-
<b>PromptFL</b>	<b>24.9<math>\pm</math>0.4</b>	58.2 $\pm$ 0.3	35.6 $\pm$ 0.6	39.6	-
<b>FedOTP</b>	23.8 $\pm$ 0.3	58.3 $\pm$ 0.8	35.2 $\pm$ 0.4	39.1	<u>1</u>
<b>FedTPG</b>	24.5 $\pm$ 0.4	58.6 $\pm$ 0.6	35.7 $\pm$ 0.3	39.6	<b>2</b>
<b><i>f</i>-CoCoOp</b>	24.0 $\pm$ 0.6	<b>59.8<math>\pm</math>1.1</b>	<b>36.0<math>\pm</math>1.0</b>	<b>39.9</b>	<b>2</b>
<b><i>f</i>-PLOT</b>	23.8 $\pm$ 0.7	57.5 $\pm$ 0.3	34.6 $\pm$ 0.6	38.6	0
<b><i>f</i>-ProDA</b>	<u>24.7<math>\pm</math>1.7</u>	58.3 $\pm$ 0.6	35.6 $\pm$ 0.8	39.5	<u>1</u>
<b><i>f</i>-ProGrad</b>	23.5 $\pm$ 1.1	58.3 $\pm$ 0.5	35.5 $\pm$ 1.2	39.1	<u>1</u>
<b><i>f</i>-SRC</b>	24.0 $\pm$ 0.7	58.7 $\pm$ 0.9	35.1 $\pm$ 0.5	39.3	1
<b><i>f</i>-KgCoOp</b>	<u>24.7<math>\pm</math>1.2</u>	<u>58.7<math>\pm</math>1.4</u>	<u>35.9<math>\pm</math>0.7</u>	<u>39.8</u>	<b>2</b>

Test-time image feature injection and self-regularization contribute to improving the robustness against cross domain shift.

## Project & Code



# Thanks

