



# Worst Case Matters for Few-Shot Recognition

Minghao Fu, Yun-Hao Cao, Jianxin Wu\*

State Key Laboratory for Novel Software Technology, Nanjing University  
fumh@lamda.nju.edu.cn, caoyh@lamda.nju.edu.cn, wujx2001@gmail.com



## 1. Introduction & Motivation

This paper focuses on improving the worst-case accuracy of few-shot learning. To accomplish this goal, we propose two strategies, i.e., stability regularization (SR) and adaptability calibration (AC) from the perspective of bias-variance tradeoff.

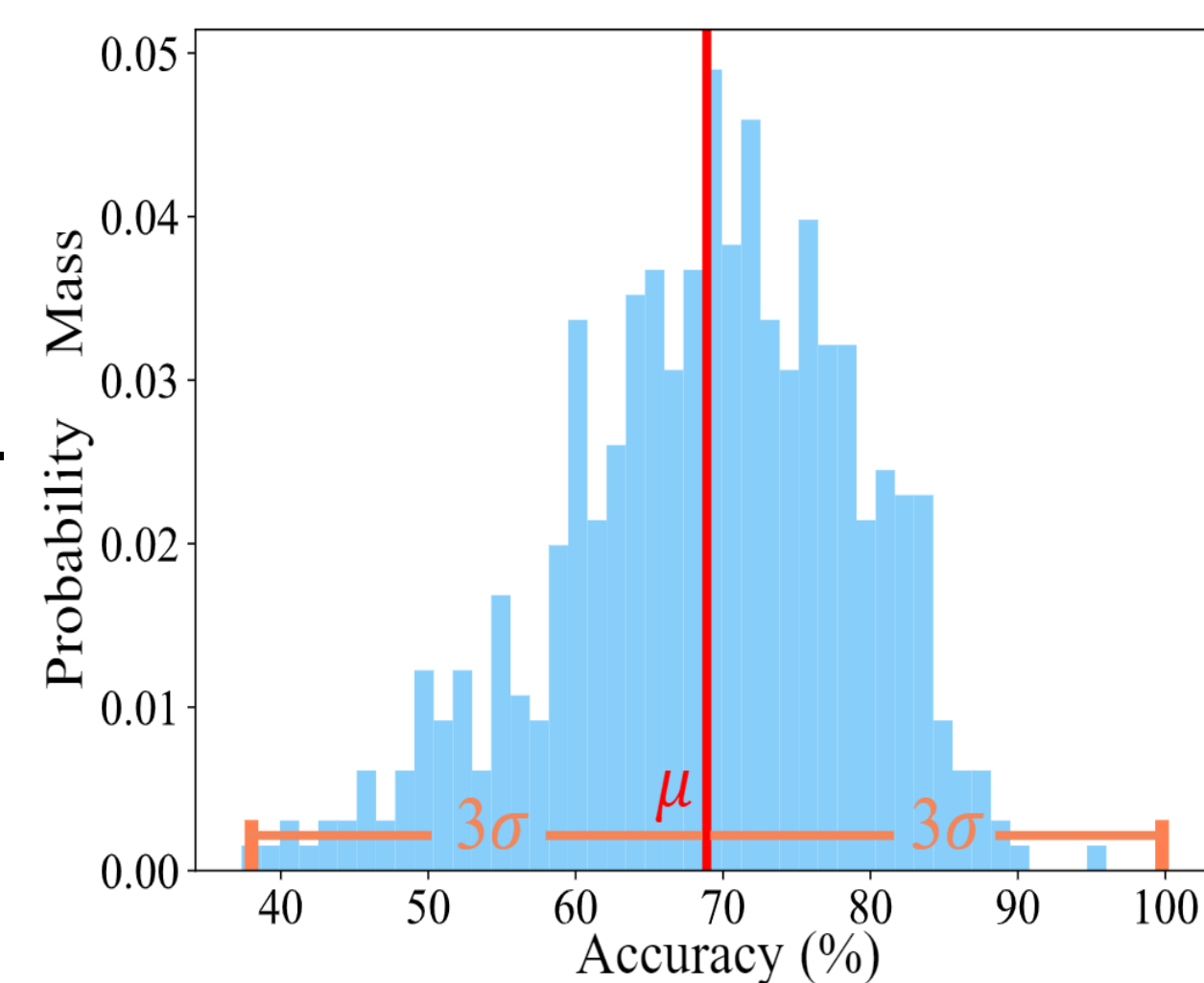
### Key idea

- We aim to boost the worst-case accuracy of few-shot learning.

### Why worst-case accuracy?

- Current criterion: Average accuracy and 95% confidence interval over many episodes.
- A high average accuracy does not necessarily mean a high worst-case accuracy.

Method	$ACC_m$	$Z_{95\%}$	$\sigma$	$ACC_1$	$ACC_{10}$
Negative-Cosine [21]	61.72	0.81	10.12	24.27	36.13
MixtFSL [2]	64.31	0.79	9.87	30.67	35.07
S2M2 <sub>R</sub> [23]	64.93	<b>0.18</b>	<b>9.18</b>	37.58	42.87
PT+NCM [16]	65.35	0.20	10.20	32.00	38.13
CGCS [11]	67.02	0.20	10.20	<b>38.70</b>	<b>44.00</b>
LR-DC [36]	<b>68.57</b>	0.55	10.28 <sup>b</sup>	37.33	42.72



- Few-shot learning is very unstable. The worst-case lags far behind the average.
- In real-world applications, worst-case accuracy is very important.

## 2. Surrogate of worst-case accuracy & Framework

- The accuracy distribution often fits well to a Gaussian. The worst-case accuracy is naturally estimated by the  $3\sigma$  rule as  $\mu - 3\sigma$ . Therefore we propose *maximizing mean accuracy*  $\mu$  and *minimizing standard deviation*  $\sigma$  **simultaneously**.
- Since directly optimizing  $\sigma$  is not plausible, according to bias-variance decomposition theory, we propose **stability regularization (SR)** with model ensemble to reduce variance, and **adaptability calibration (AC)** to reduce bias.

### ✓ Stability Regularization:

- A loss function to prevent model from completely forgetting the representation learned from the base set.

$$\mathcal{L}_S(\mathbf{x}) = -\frac{f(\mathbf{x}) \cdot \hat{f}(\mathbf{x})}{\|f(\mathbf{x})\| \|\hat{f}(\mathbf{x})\|}$$

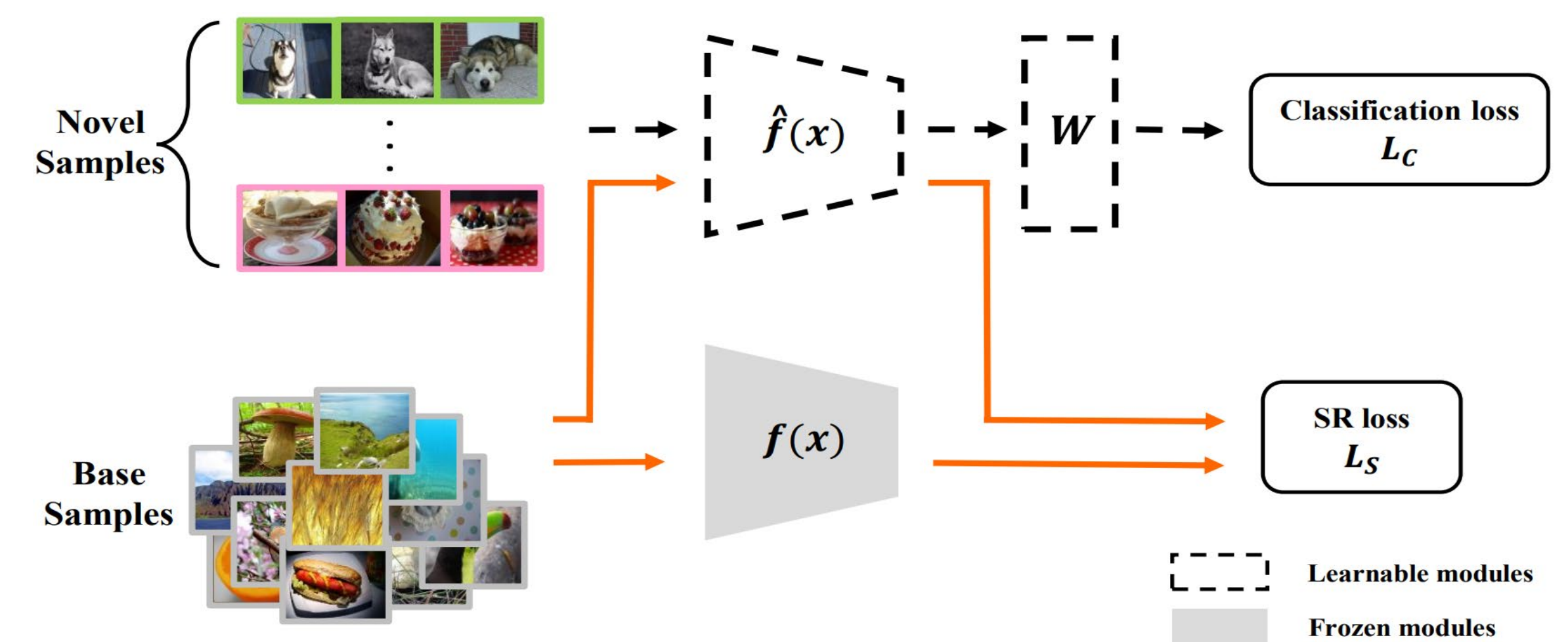
- Generalizable to other unlabeled images.
- Ensemble.

### ✓ Adaptability Calibration:

- A remedy to properly calibrated capacity increase.

### Bias-variance Tradeoff of SR and AC

- SR: More stable classifier has a smaller variance ( $\sigma$ ).
- AC: Larger model capacity has a smaller bias. Because of the contradiction of two terms, a balance is necessary.



## 3. Experiments

### State-of-the-art comparison

Dataset	Method	1-shot			5-shot			
		$ACC_1$	$\mu - 3\sigma$	$\sigma$	$ACC_m$	$ACC_1$	$\mu - 3\sigma$	$\sigma$
mini-ImageNet	ProtoNet <sup>†</sup> [29]	19.76	23.41	10.25	54.16	43.74	49.32	8.12
	Transductive-FT [7]	24.00	32.82	10.97	65.73	50.67	53.23	8.39
	Negative-Cosine [21]	24.27	31.36	10.12	61.72	53.30	61.18	6.87
	MixtFSL [2]	30.67	34.70	9.87	64.31	46.67	59.16	7.50
	PT+NCM [16]	32.00	34.75	10.20	65.35	56.00	63.98	6.63
	LR-DC [36]	37.33	37.73	10.28	68.57	60.52	63.02	6.62
	S2M2 <sub>R</sub> [23]	37.58	37.39	<b>9.18</b>	64.93	58.66	66.35	<b>5.61</b>
	CGCS [11]	38.70	36.42	10.20	67.02	49.30	60.90	7.14
	AC+SR (ours)	40.52	40.25	9.71	69.38	63.20	66.46	6.47
	AC+EnSR (ours)	<b>40.52</b>	<b>40.67</b>	9.64	<b>69.59</b>	<b>63.48</b>	<b>66.71</b>	6.42
CUB	ProtoNet <sup>†</sup> [29]	28.00	39.99	11.00	72.99	53.33	67.53	6.37
	Negative-Cosine [21]	36.00	40.80	10.62	72.66	70.70	73.29	5.37
	PT+NCM [16]	40.00	49.97	10.20	80.57	69.33	75.85	5.10
	MixtFSL [2]	40.00	32.69	13.75	73.94	57.33	67.26	6.25
	LR-DC [36]	44.00	49.74	9.94	79.56	68.80	75.49	5.06
	CGCS [11]	50.67	42.53	10.71	74.66	57.33	70.01	6.12
	S2M2 <sub>R</sub> [23]	52.00	50.32	10.12	80.68	73.86	74.35	5.50
	AC+SR (ours)	52.78	58.44	8.90	85.14	76.00	81.85	4.19
	AC+EnSR (ours)	<b>53.04</b>	<b>58.84</b>	<b>8.86</b>	<b>85.42</b>	<b>77.58</b>	<b>82.14</b>	<b>4.13</b>

- Our method performs best on worst-case and average accuracy.
- $\mu - 3\sigma$  correlates better with  $ACC_1$  than  $\sigma$  or  $\mu$ .

### Results for different degree of AC

$W$	AC				1-shot			5-shot		
	5	4	3	2+1	$ACC_m$	$\sigma$	$ACC_1$	$ACC_{10}$	$ACC_{100}$	$ACC_m$
✓					59.53	<b>9.99</b>	30.42	35.09	45.05	73.61
✓✓					62.76	10.16	<b>32.78</b>	<b>37.55</b>	<b>47.84</b>	79.23
✓✓✓					62.77	10.38	30.12	35.44	47.46	79.91
✓✓✓✓					<b>62.80</b>	10.50	22.96	34.93	47.27	<b>80.04</b>
✓✓✓✓✓					62.19	10.50	25.60	35.29	46.68	79.87

- Fine-tuning  $W+$  'res5' is the AC strategy used in our method.

### Generalize SR to other data

$D_{sr}$	$ACC_m$	$\sigma$	$ACC_1$	$ACC_{10}$	$ACC_{100}$
-	60.09	10.13	25.60	34.55	45.24
mini-ImageNet	<b>62.33</b>	10.27	32.00	<b>37.03</b>	<b>47.36</b>
CUB [34]	60.66	9.99	30.14	36.43	46.29
Cars [18]	61.05	10.29	30.94	36.64	46.20
DTD [6]	61.99	10.23	31.20	36.75	47.23
Pets [25]	61.57	<b>9.95</b>	32.00	36.62	47.27
Flower [24]	61.03	10.07	<b>32.26</b>	36.51	46.37
CIFAR-100 [19]	60.47	10.04	30.68	36.69	46.02

- Don't require images used in SR to be visually similar or semantically correlated to support samples.
- SR is consistently useful.

## 4. Contributions & Conclusions

- ✓ We are the first to emphasize the importance and to advocate the adoption of **worst case accuracy** in few-shot learning.
- ✓ We argue that in addition to maximizing the average accuracy  $\mu$ , we must also **simultaneously** reduce the standard deviation  $\sigma$ .
- ✓ We propose stability regularization (SR) with model ensemble and adaptability calibration (AC) strategies from the perspective of bias-variance tradeoff.
- ✓ Our method shows superior results compared to current state-of-the-art methods in terms of not only average, but also worst-case accuracy.