

# Meta Dropout

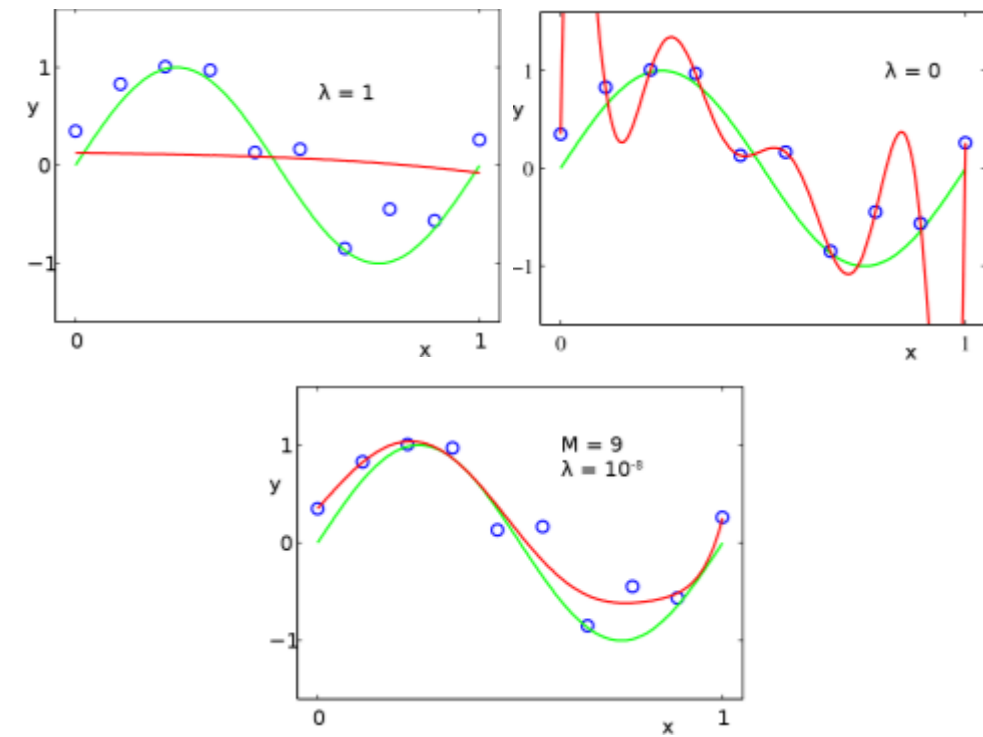
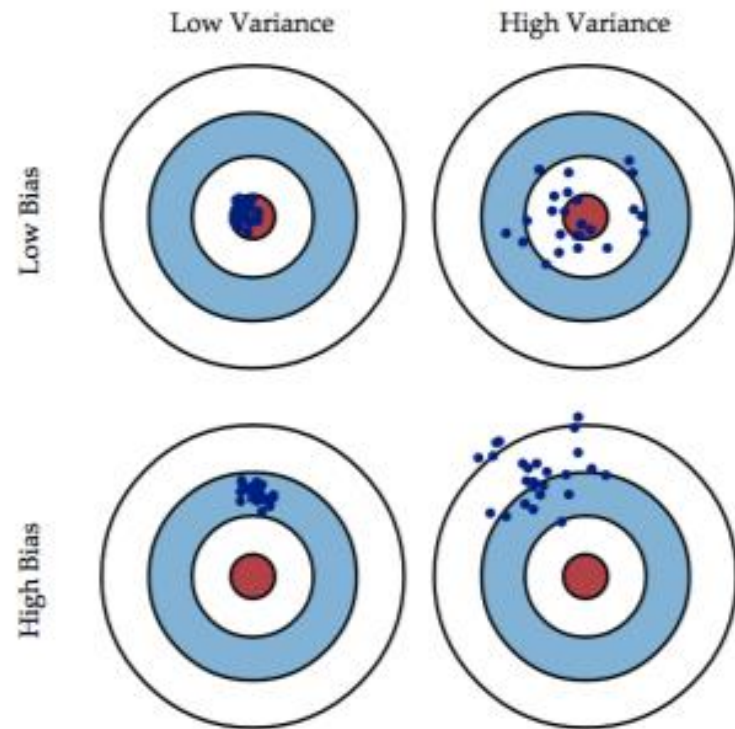
# Learning to Perturb Features for Generalization

Hae Beom Lee, Taewook Nam, Eunho Yang, Sung Ju Hwang

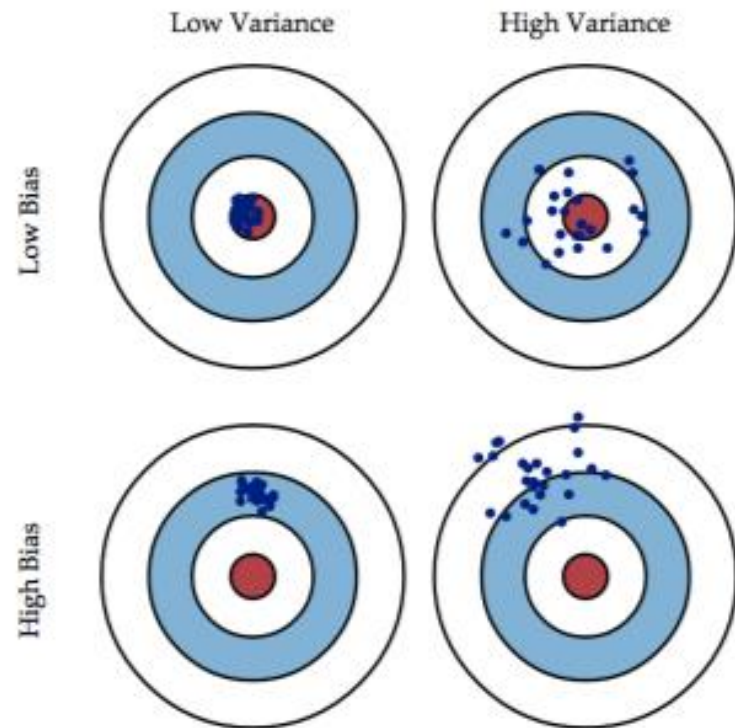
Tobias Schmidt

Munich, 7<sup>th</sup> July 2021

# Motivation – Generalization



# Motivation – Generalization



## Common Solutions - Bias/Variance Trade-Off

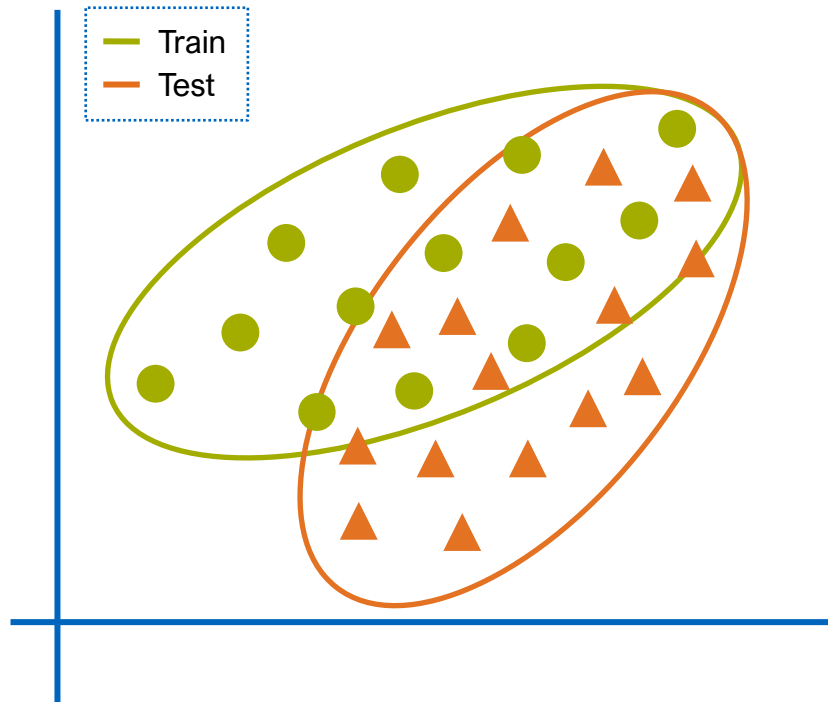
### Appropriate Priors

- Image Convolutions (Shift-Invariance)
- Graph Convolutions (Permutation-Invariance)
- ...

### Regularization

- Reducing model capacity
- Reducing information from inputs
- Smoothing loss surface
- Multi-task training
- Meta-Learning

# Motivation – Train/Test Distribution Mismatch



## Common Solutions - Bias/Variance Trade-Off

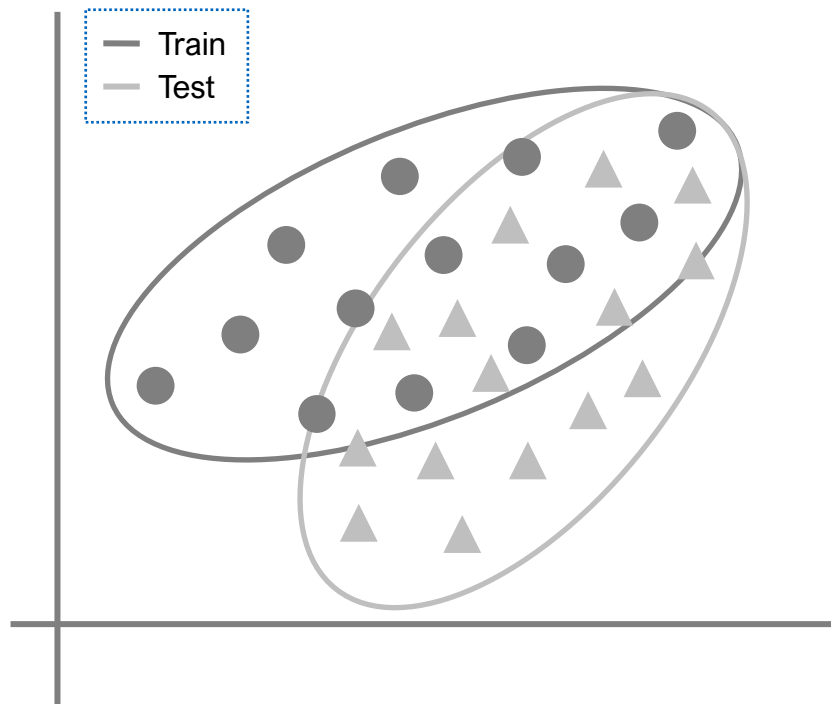
### Appropriate Priors

- Image Convolutions (Shift-Invariance)
- Graph Convolutions (Permutation-Invariance)
- ...

### Regularization

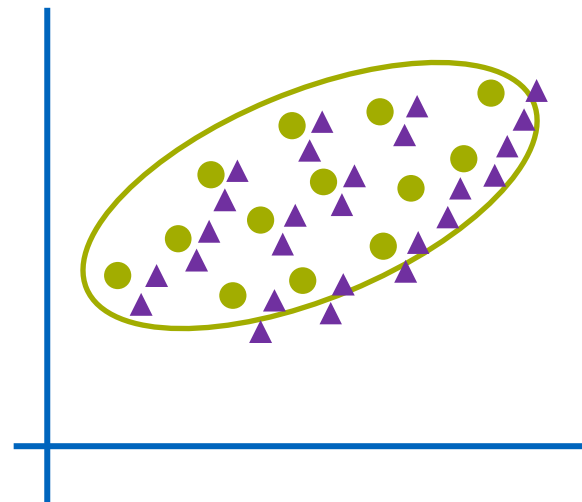
- Reducing model capacity
- Reducing information from inputs
- Smoothing loss surface
- Multi-task training
- Meta-Learning

# Motivation – Train/Test Distribution Mismatch

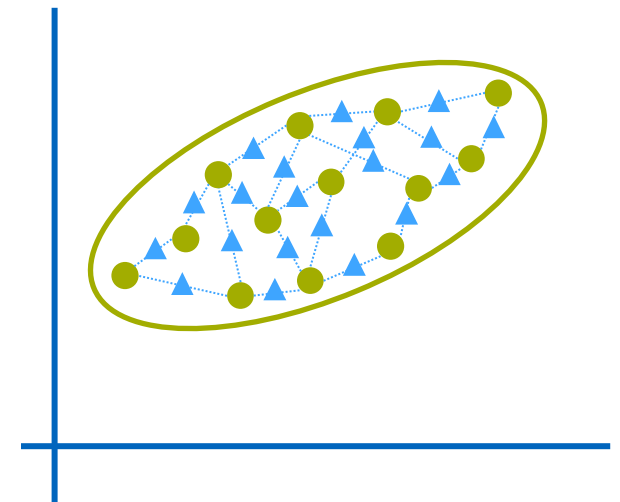


## Alternative Solutions - Simulate Test Samples

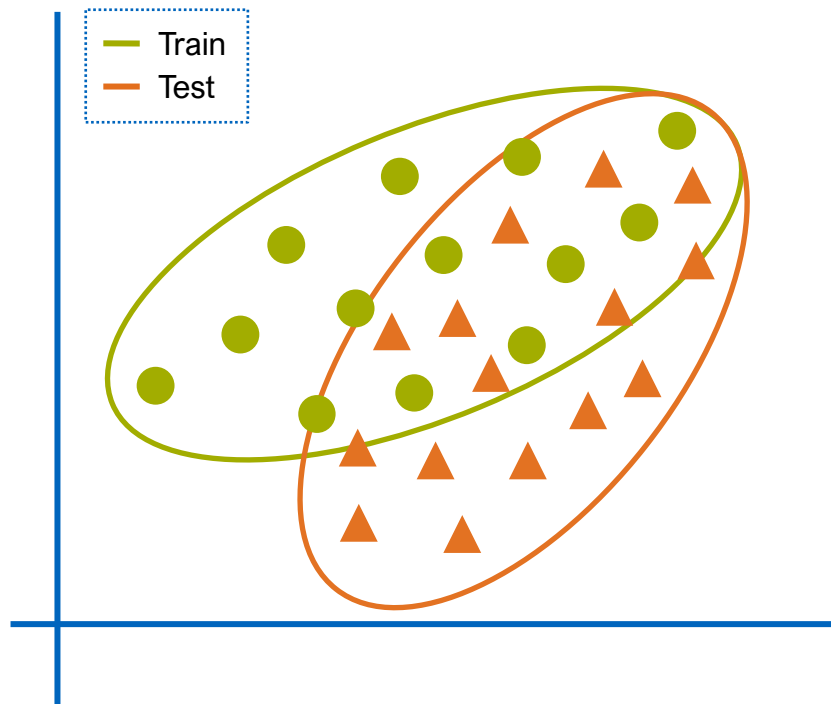
Data augmentations



Mixup

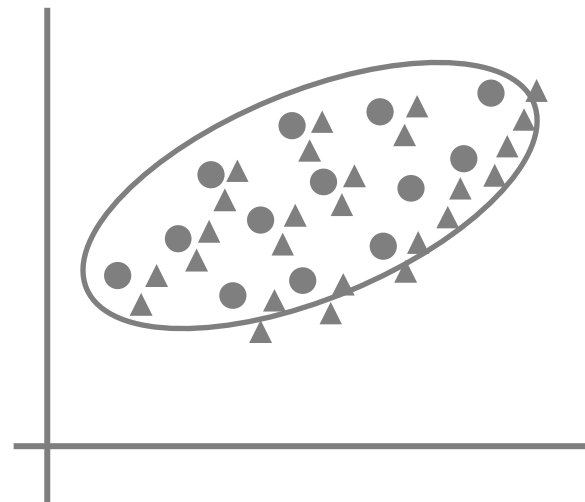


# Motivation – Train/Test Distribution Mismatch

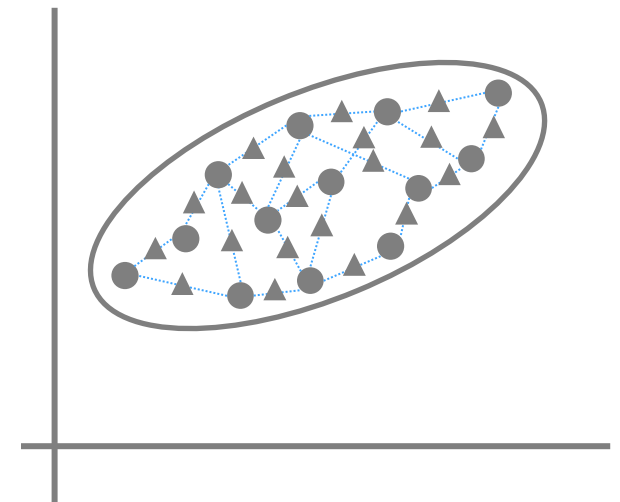


## Alternative Solutions - Simulate Test Samples

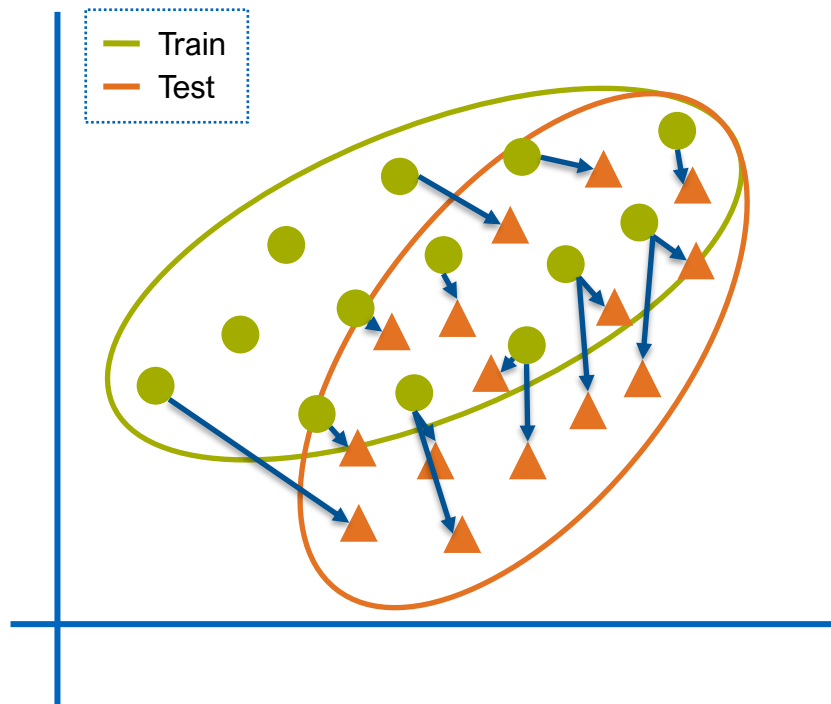
Data augmentations



Mixup

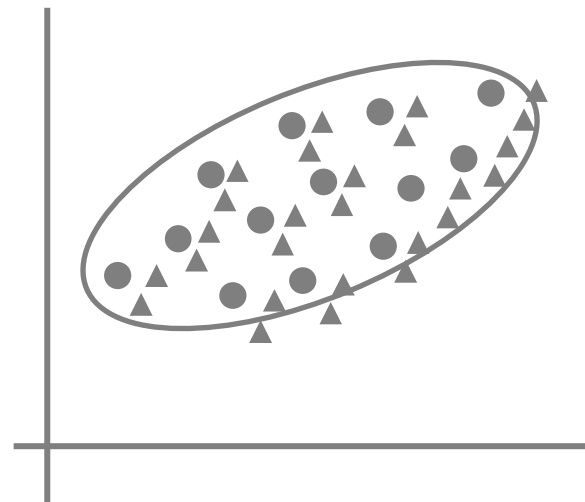


# Motivation – Train/Test Distribution Mismatch

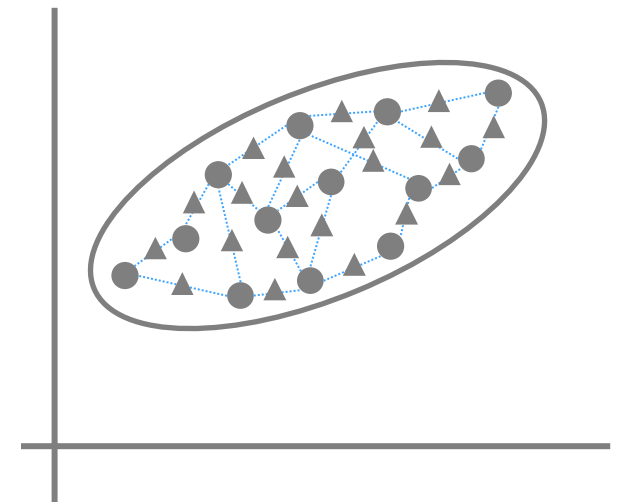


## Alternative Solutions - Simulate Test Samples

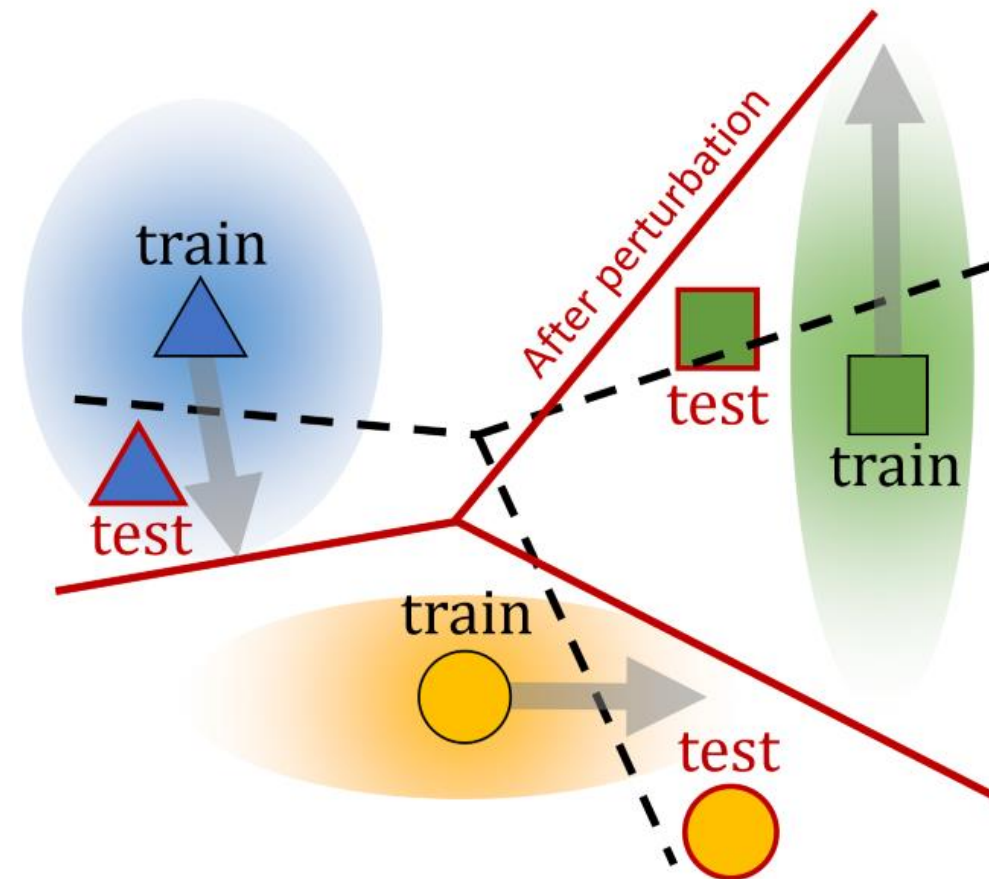
Data augmentations



Mixup



Idea: Learn to perturb the data for better generalization

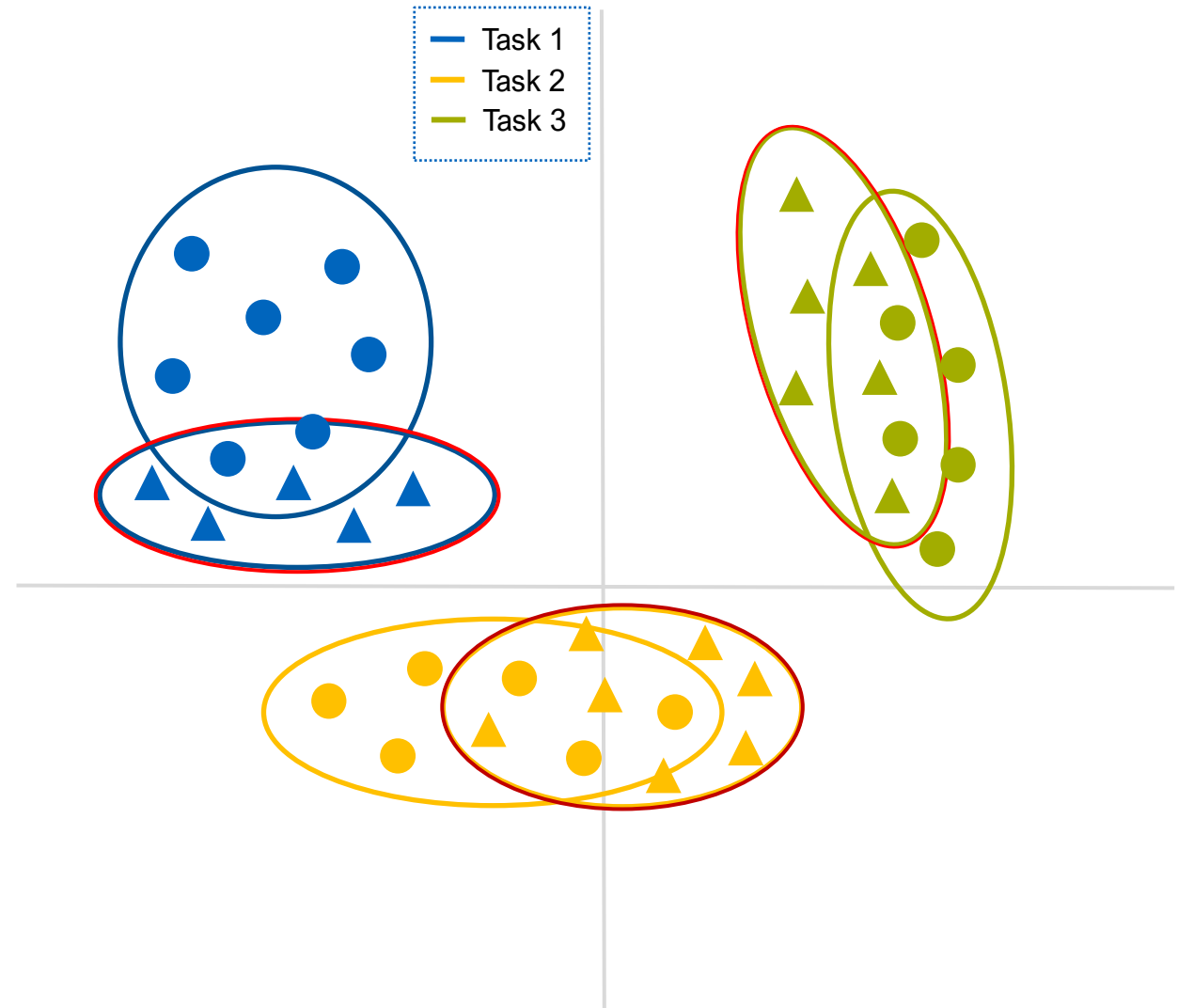




# Challenges

A training instance may need to cover multiple test instances

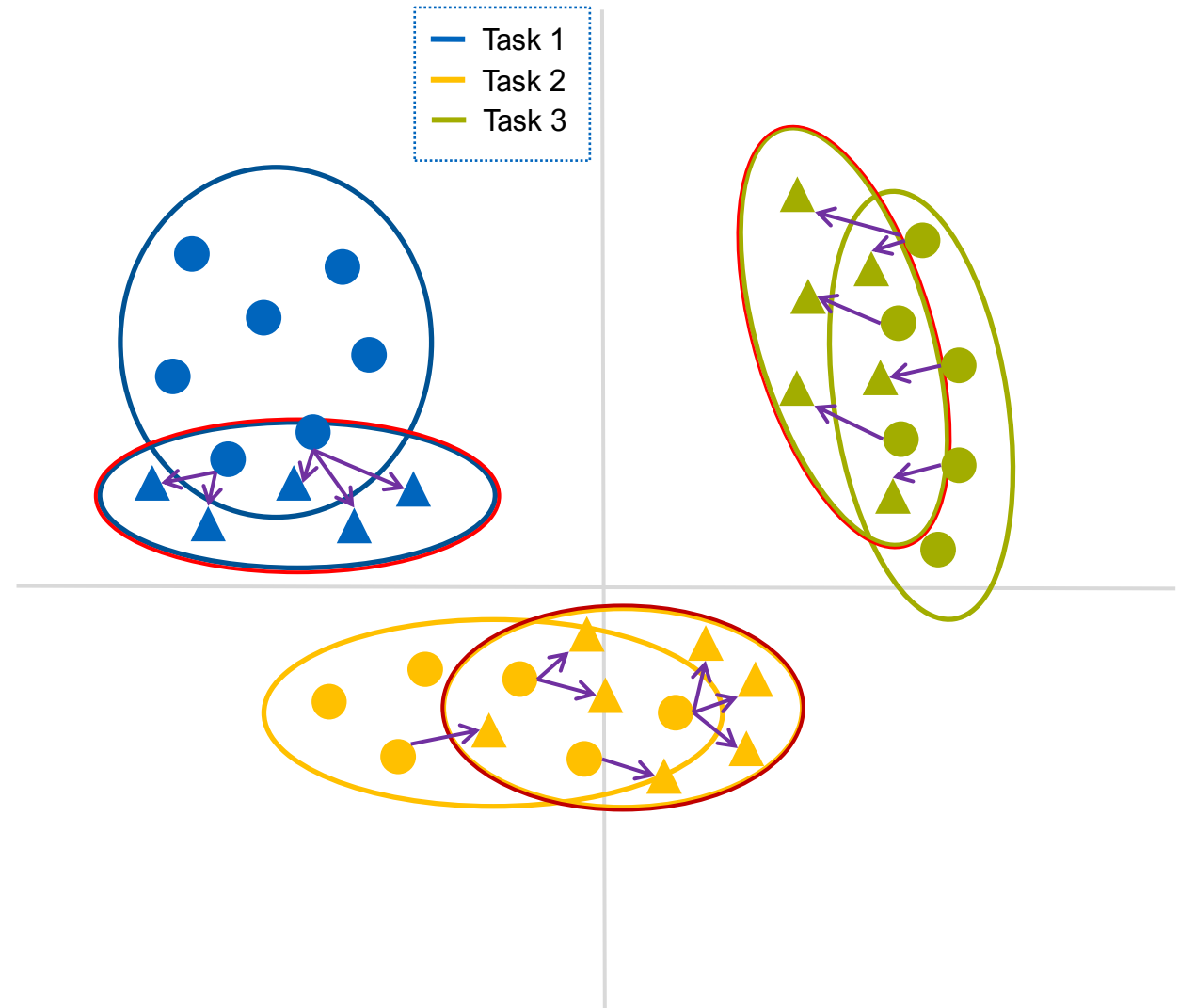
Meaningful directions differ from one task to another



# Challenges

**A training instance may need to cover multiple test instances**

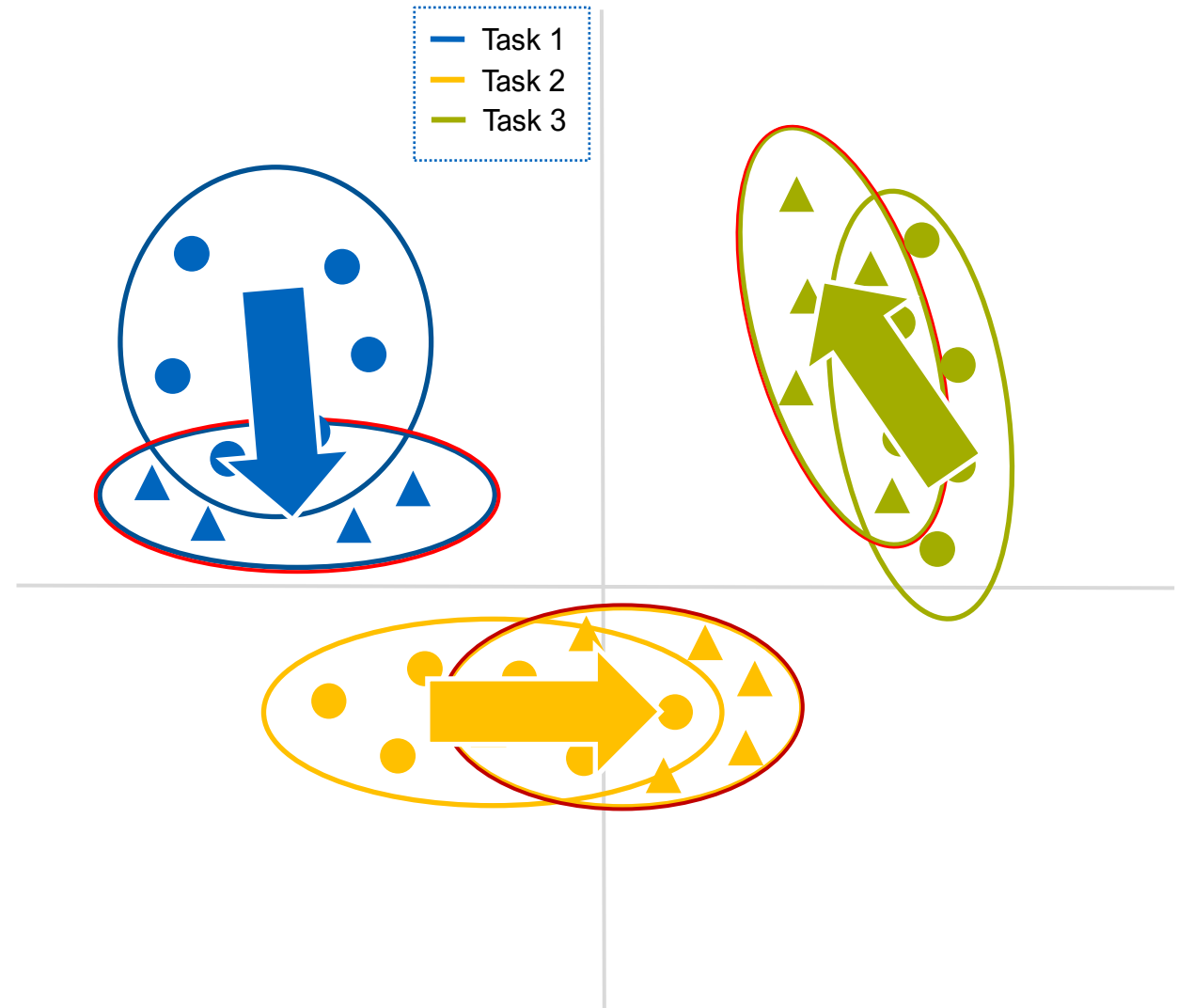
Meaningful directions differ from one task to another



# Challenges

A training instance may need to cover multiple test instances

Meaningful directions differ from one task to another



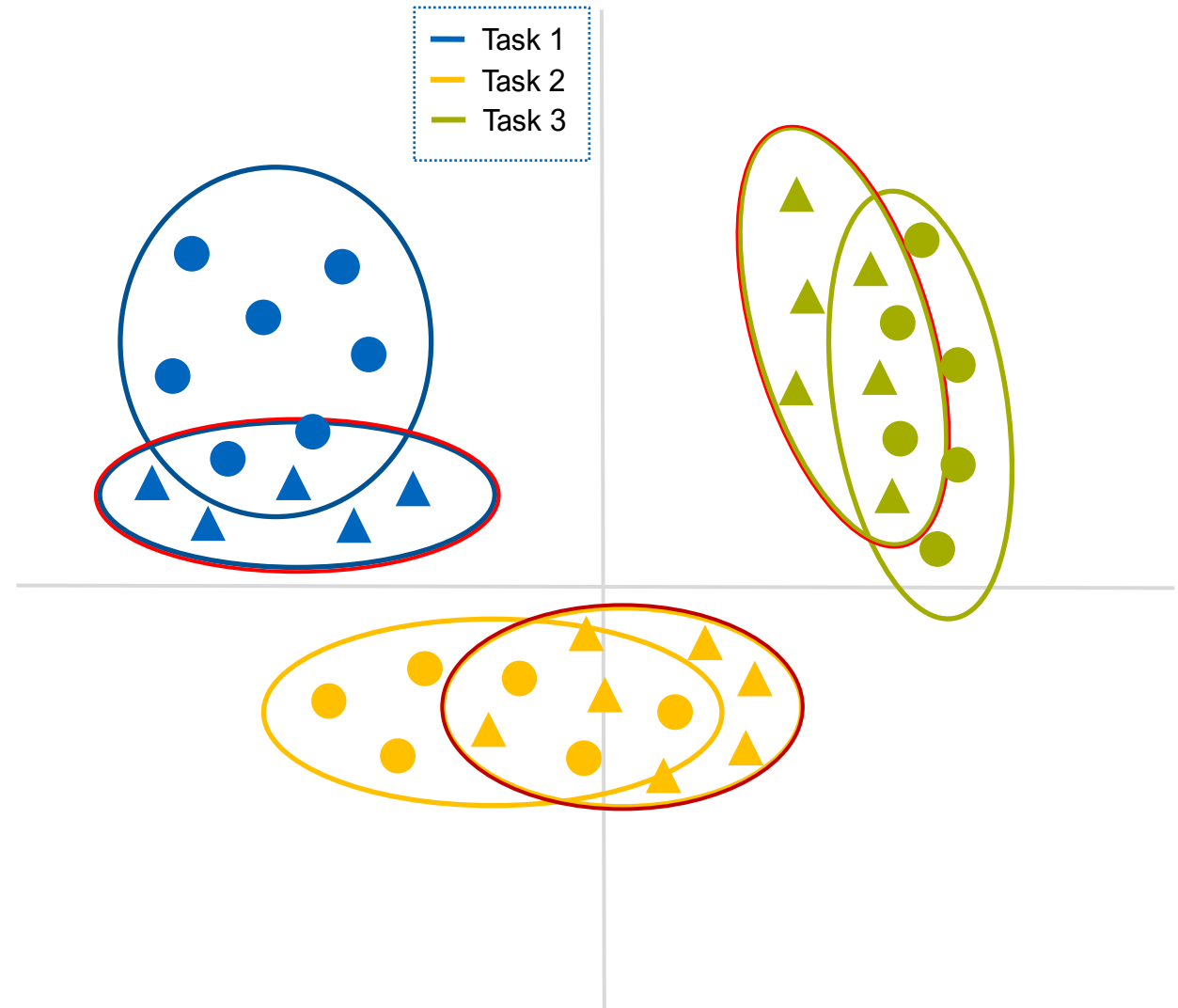
# Challenges

**A training instance may need to cover multiple test instances**

→ Noise Distribution

**Meaningful directions differ from one task to another**

→ Input-Dependent Noise



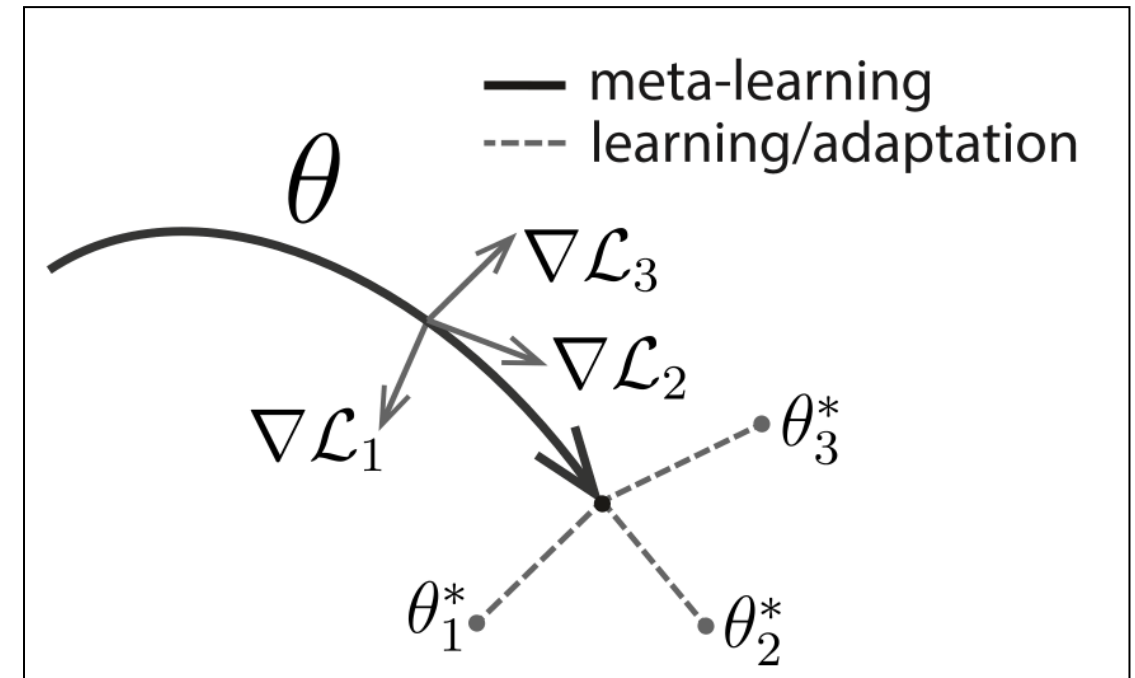
# Meta Learning Framework - MAML

## Properties / Limitations of MAML

knowledge transfer via learned parameter  $\theta$

parameter  $\theta$  only implicitly captures test distributions

→ Misses out on important knowledge about task distribution



# Recap: MAML

---

## Algorithm 1 Model-Agnostic Meta-Learning

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

1: randomly initialize  $\theta$

2: **while** not done **do**

3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$

4:   **for all**  $\mathcal{T}_i$  **do**

5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples

6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$

7:   **end for**

8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$

9: **end while**

---

Update over task  
distribution:

- **initial parameter**  $\theta$

Outer Loop

Inner Loop

Perform Few-Shot  
Learning for each task

# Meta Dropout

---

## Algorithm 1 Meta-training

---

```

1: Input: Task distribution  $p(\mathcal{T})$ , Number of inner steps  $K$ ,
2: while not converged do
3:   Sample  $(\mathcal{D}^{\text{tr}}, \mathcal{D}^{\text{te}}) \sim p(\mathcal{T})$ 
4:    $\theta_0 \leftarrow \theta$ 
5:   for  $k = 0$  to  $K - 1$  do
6:     Sample  $\tilde{z}_i \sim p(z_i | x_i^{\text{tr}}; \phi, \theta_k)$  for  $i = 1, \dots, N$ 
7:      $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta_k} \frac{1}{N} \sum_{i=1}^N \log p(y_i^{\text{tr}} | x_i^{\text{tr}}, \tilde{z}_i; \theta_k)$ 
8:   end for
9:    $\theta^* \leftarrow \theta_K$ 
10:   $\theta \leftarrow \theta + \beta \frac{1}{M} \sum_{j=1}^M \nabla_{\theta} \log p(y_j^{\text{te}} | x_j^{\text{te}}, z_j = \bar{z}_j; \theta^*)$ 
11:   $\phi \leftarrow \phi + \beta \frac{1}{M} \sum_{j=1}^M \nabla_{\phi} \log p(y_j^{\text{te}} | x_j^{\text{te}}, z_j = \bar{z}_j; \theta^*)$ 
12: end while

```

---

Update over task  
distribution:

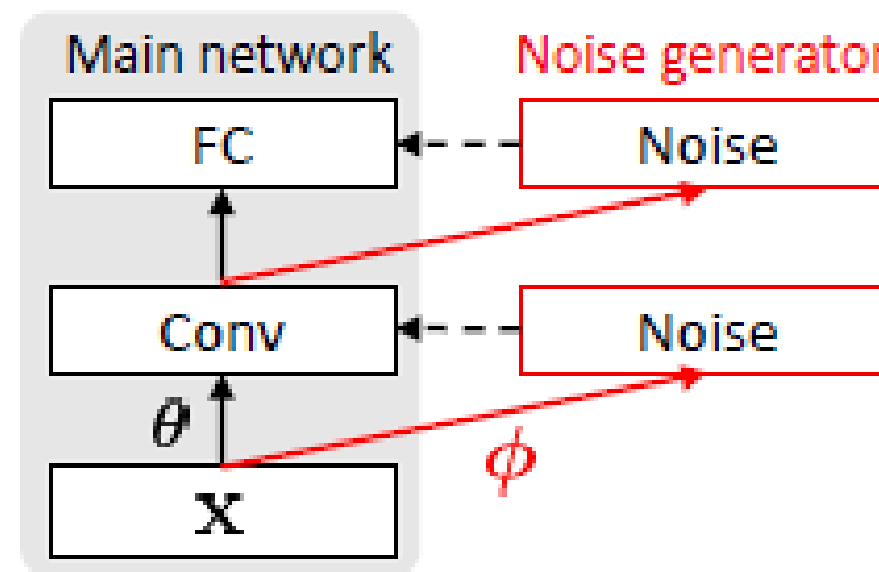
- initial parameter  $\theta$
- parameters  $\phi$  of noise  $p(z)$

Outer Loop

Inner Loop

Perform Few-Shot  
Learning for each task  
**perturbing the input  $x_i^{\text{tr}}$   
with multiplicative  
noise  $\tilde{z}_i$**

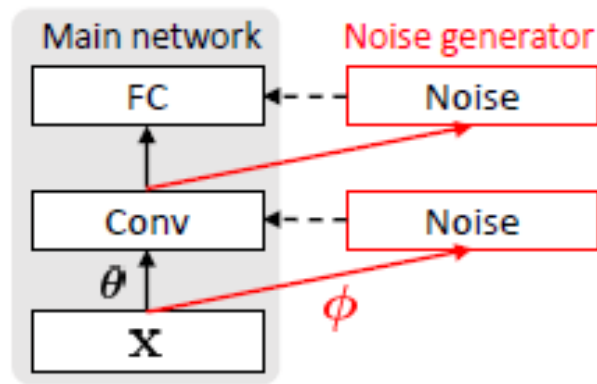
# Meta-Dropout: Model Architecture





# Meta-Dropout: Implementation Detail

## Model Architecture



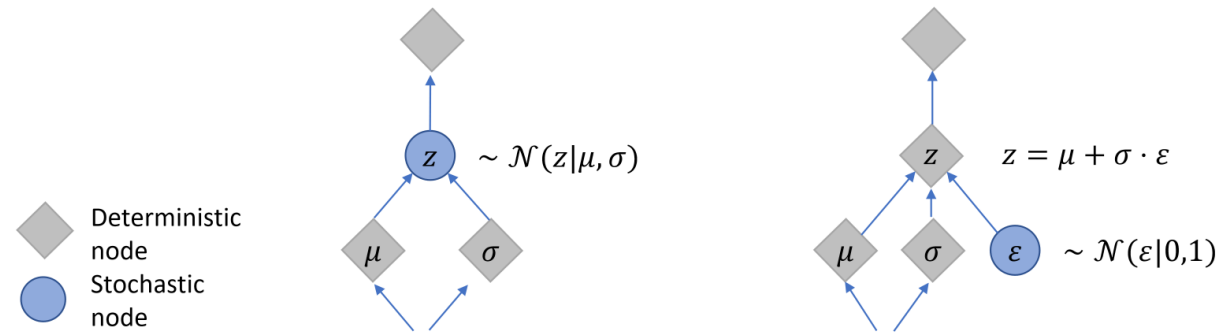
## Reparameterization Trick

$$\log p(Y_i^{tr} | X_i^{tr}, \theta, \phi) \geq \sum_{i=1}^N \mathbb{E}_{z_i \sim p(z_i | x_i^{tr}, \theta, \phi)} [\log p(y_i^{tr} | x_i^{tr}, \theta, \phi)]$$

$$\approx \sum_{i=1}^N \sum_{s=1}^S \log p(y_i^{tr} | x_i^{tr}, z_i^{(s)}, \theta) \text{ with } z_i^{(s)} \sim p(z_i | x_i^{tr}, \theta, \phi)$$

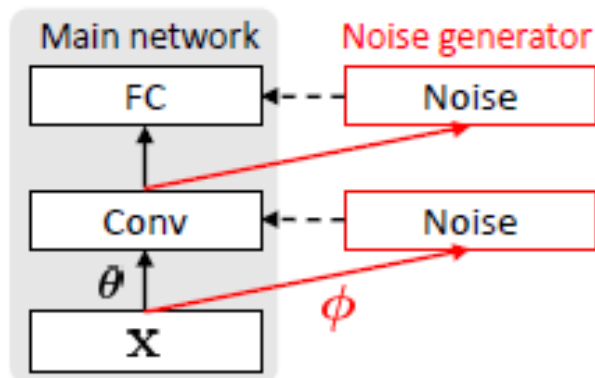
Approximation via Monte-Carlo

Impossible to calculate gradient!



# Meta-Dropout: Implementation Detail

## Model Architecture



**Form of the Noise:**  $p(y_i^{tr} | x_i^{tr}, z_i^{(s)}, \theta)$  with  $z_i^{(s)} \sim p(z_i | x_i^{tr}, \theta, \phi)$

### Additive Noise

$$\begin{aligned} h^{(0)} &= x_i^{tr} \\ h^{(l)} &= \text{ReLU}(f^{(l)}(h^{(l-1)}) + z^{(l)}) \end{aligned} \quad z_i^{(s)} \sim N(z^{(l)} | 0, \lambda^2 \text{diag}(\sigma^2))$$

### Multiplicative Noise

$$\begin{aligned} h^{(0)} &= x_i^{tr} \\ h^{(l)} &= \text{ReLU}(f^{(l)}(h^{(l-1)}) \circ z^{(l)}) \end{aligned} \quad \begin{aligned} z^{(l)} &= \text{Softplus}(a^{(l)}) \\ a^{(l)} &\sim N(a^{(l)} | \mu^{(l)}, I) \end{aligned}$$

# Experiments - Datasets

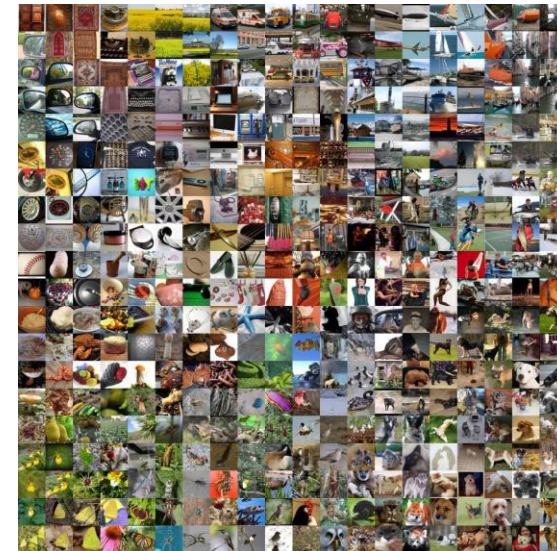
## Omniglot

Handwritten character classification  
→ 20 instances of ~1600 characters from 50 alphabets



## minImageNET

Small version of ImageNET  
→ 100 classes with 600 samples



# Experiments

## Meta-Learning Frameworks

MAML

Meta-SGD

Prototypical Networks

Matching Networks

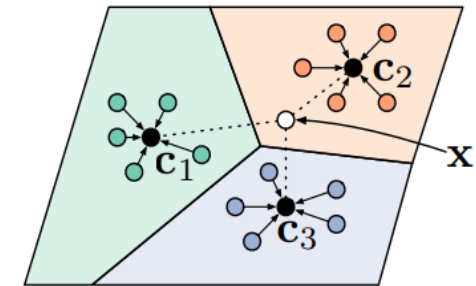
Reptile

Amortized Bayesian ML

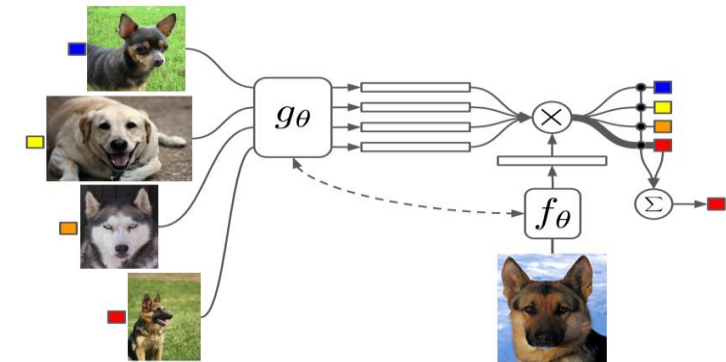
Probabilistic MAML

MT-NET

CAVIA



Prototypical Networks



Matching Networks

# Experiments

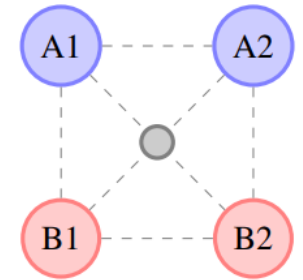
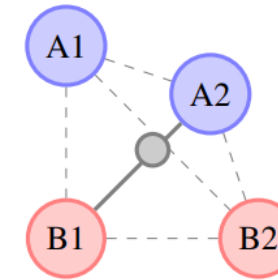
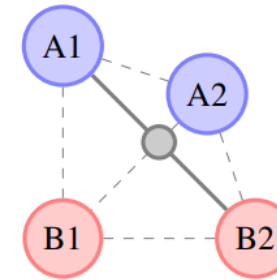
## Perturbation Based Methods:

Input & Manifold Mixup

Variational Information Bottleneck

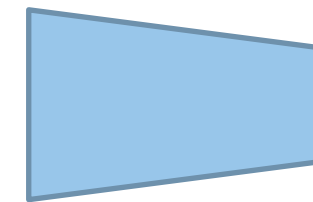
Information Dropout

Adversarial learning/Training

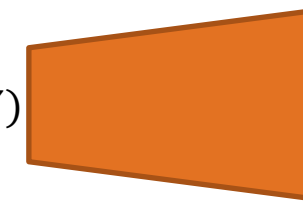


Input & Manifold Mixup

$X$



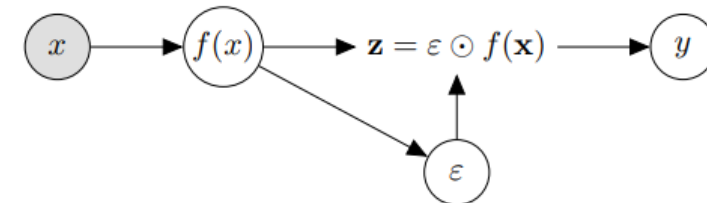
$p(Z)$



$Y$

Encode maximal information about Target  $Y$  in latent stochastic encoding  $Z$  measured by mutual information  $I(Z, Y) < I_c$  where  $I_c$  is information constraint

Variational Information Bottleneck



Information Dropout

# Few-shot classification performance

Models	Omniglot 20-way		miniImageNet 5-way	
	1-shot	5-shot	1-shot	5-shot
Meta-Learning LSTM (Ravi & Larochelle, 2017)	-	-	43.44 $\pm$ 0.77	60.60 $\pm$ 0.71
Matching Networks (Vinyals et al., 2016)	93.8	98.7	43.56 $\pm$ 0.84	55.31 $\pm$ 0.73
Prototypical Networks (Snell et al., 2017)	95.4	98.7	46.14 $\pm$ 0.77	65.77 $\pm$ 0.70
Prototypical Networks (Snell et al., 2017) (Higher way)	96.0	98.9	49.42 $\pm$ 0.78	<b>68.20<math>\pm</math>0.66</b>
MAML (our reproduction)	95.23 $\pm$ 0.17	98.38 $\pm$ 0.07	49.58 $\pm$ 0.65	64.55 $\pm$ 0.52
Meta-SGD (our reproduction)	96.16 $\pm$ 0.14	98.54 $\pm$ 0.07	48.30 $\pm$ 0.64	65.55 $\pm$ 0.56
Reptile (Nichol et al., 2018)	89.43 $\pm$ 0.14	97.12 $\pm$ 0.32	49.97 $\pm$ 0.32	65.99 $\pm$ 0.58
Amortized Bayesian ML (Ravi & Beatson, 2019)	-	-	45.00 $\pm$ 0.60	-
Probabilistic MAML (Finn et al., 2018)	-	-	50.13 $\pm$ 1.86	-
MT-Net (Lee & Choi, 2018)	96.2 $\pm$ 0.4	-	51.70 $\pm$ 1.84	-
CAVIA (512) (Zintgraf et al., 2019)	-	-	51.82 $\pm$ 0.65	65.85 $\pm$ 0.55
<b>MAML + Meta-dropout</b>	<b>96.63<math>\pm</math>0.13</b>	<b>98.73<math>\pm</math>0.06</b>	<b>51.93<math>\pm</math>0.67</b>	<b>67.42<math>\pm</math>0.52</b>
<b>Meta-SGD + Meta-dropout</b>	<b>97.02<math>\pm</math>0.13</b>	<b>99.05<math>\pm</math>0.05</b>	50.87 $\pm$ 0.63	65.55 $\pm$ 0.57

## Few-shot classification performance

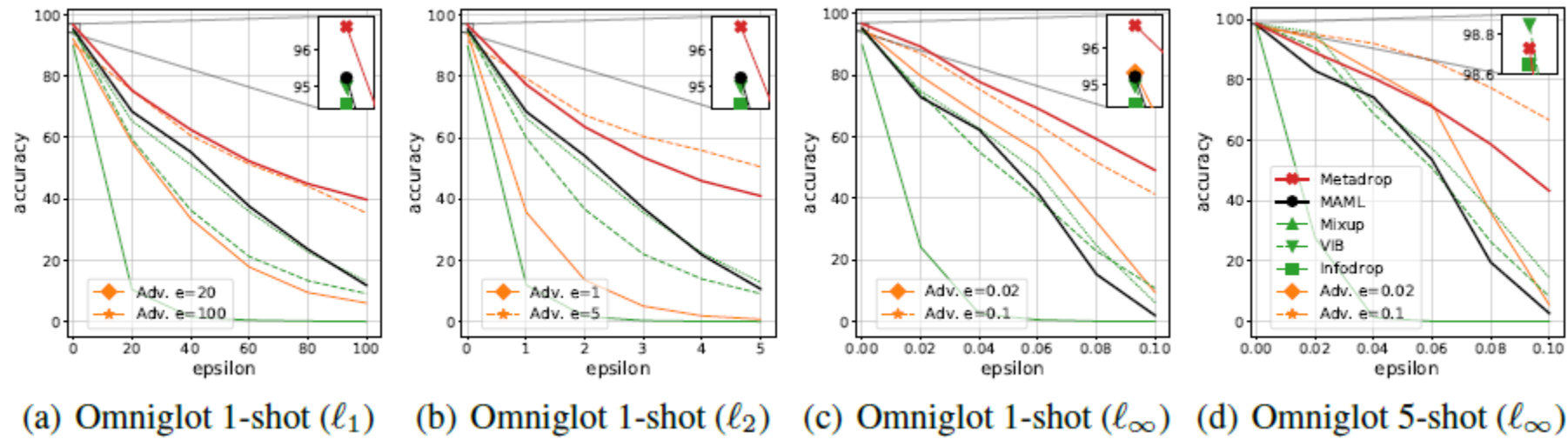
Models (MAML +)	Noise Type	Hyper- parameter	Omniglot 20-way		miniImageNet 5-way	
			1-shot	5-shot	1-shot	5-shot
No perturbation		None	95.23 $\pm$ 0.17	98.38 $\pm$ 0.07	49.58 $\pm$ 0.65	64.55 $\pm$ 0.52
Input & Manifold Mixup (Zhang et al., 2017)	Pairwise	$\gamma = 0.2$	89.78 $\pm$ 0.25	97.86 $\pm$ 0.08	48.62 $\pm$ 0.66	63.86 $\pm$ 0.53
(Verma et al., 2019)		$\gamma = 1$	87.00 $\pm$ 0.28	97.27 $\pm$ 0.10	48.24 $\pm$ 0.62	62.32 $\pm$ 0.54
		$\gamma = 2$	87.26 $\pm$ 0.28	97.14 $\pm$ 0.17	48.42 $\pm$ 0.64	62.56 $\pm$ 0.55
Variational	Add.	$\beta = 10^{-5}$	92.09 $\pm$ 0.22	98.85 $\pm$ 0.07	48.12 $\pm$ 0.65	64.78 $\pm$ 0.54
Information Bottleneck		$\beta = 10^{-4}$	93.01 $\pm$ 0.20	98.80 $\pm$ 0.07	46.75 $\pm$ 0.63	64.07 $\pm$ 0.54
(Alemi et al., 2017)		$\beta = 10^{-3}$	94.98 $\pm$ 0.16	98.75 $\pm$ 0.07	47.59 $\pm$ 0.60	63.30 $\pm$ 0.53
Information Dropout	Mult.	$\beta = 10^{-5}$	94.49 $\pm$ 0.17	98.50 $\pm$ 0.07	50.36 $\pm$ 0.68	65.91 $\pm$ 0.55
(ReLU ver.)		$\beta = 10^{-4}$	94.36 $\pm$ 0.17	98.53 $\pm$ 0.07	49.14 $\pm$ 0.63	64.96 $\pm$ 0.54
(Achille & Soatto, 2018)		$\beta = 10^{-3}$	94.28 $\pm$ 0.17	98.65 $\pm$ 0.07	43.78 $\pm$ 0.61	63.36 $\pm$ 0.56
Meta-dropout	Add.	0.1	96.55 $\pm$ 0.14	99.04 $\pm$ 0.05	50.25 $\pm$ 0.66	66.78 $\pm$ 0.53
(See Appendix B for Add.)	Mult.	None	96.63 $\pm$ 0.13	98.73 $\pm$ 0.06	51.93 $\pm$ 0.67	67.42 $\pm$ 0.52

# Ablation study on the noise type

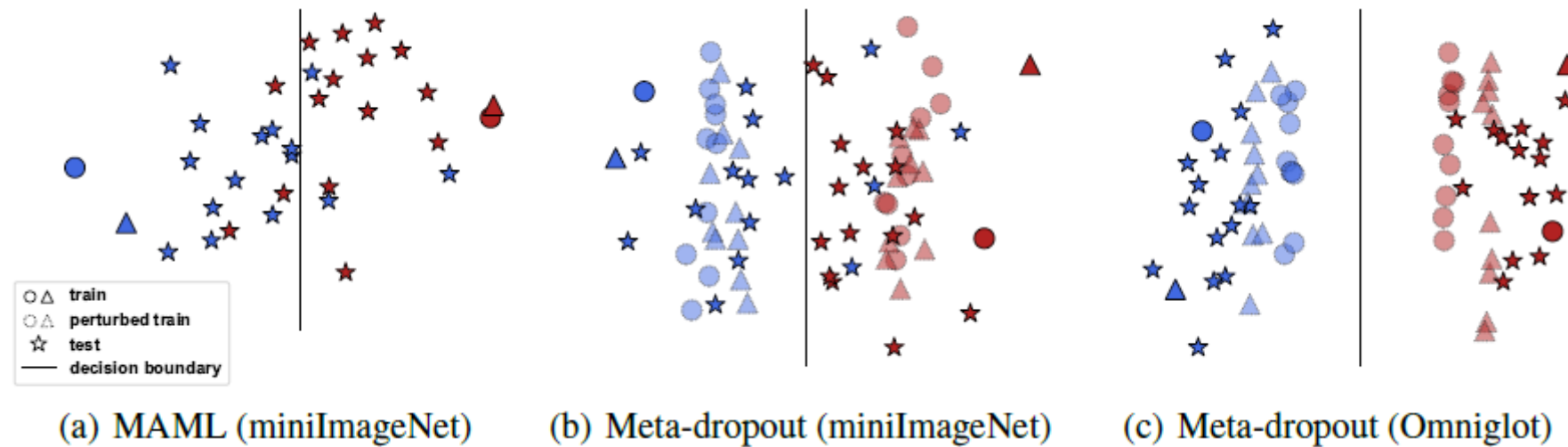
Models (MAML+)	Rand. samp.	Learn. mult.	Input dep.	Omniglot 20-way 1-shot	Omniglot 20-way 5-shot	miniImageNet 5-way 1-shot	miniImageNet 5-way 5-shot
None	X	X	X	95.23 $\pm$ 0.17	98.38 $\pm$ 0.07	49.58 $\pm$ 0.65	64.55 $\pm$ 0.52
Fixed Gaussian (✓)	O	X	X	95.44 $\pm$ 0.17	<b>98.99<math>\pm</math>0.06</b>	49.39 $\pm$ 0.63	66.84 $\pm$ 0.54
Weight Gaussian	O	X	X	94.32 $\pm$ 0.18	98.35 $\pm$ 0.07	49.37 $\pm$ 0.64	64.78 $\pm$ 0.54
Independent Gaussian	O	O	X	94.36 $\pm$ 0.18	98.26 $\pm$ 0.08	50.31 $\pm$ 0.64	66.97 $\pm$ 0.54
MAML + More param	X	O	O	95.83 $\pm$ 0.15	97.85 $\pm$ 0.09	50.63 $\pm$ 0.64	65.20 $\pm$ 0.51
Determ. Meta-drop. (✓)	X	O	O	95.99 $\pm$ 0.14	97.78 $\pm$ 0.09	50.75 $\pm$ 0.63	65.62 $\pm$ 0.53
Meta-drop. w/ learned var.	O	O	O	95.98 $\pm$ 0.15	98.87 $\pm$ 0.06	50.93 $\pm$ 0.68	66.15 $\pm$ 0.56
<b>Meta-dropout</b>	O	O	O	<b>96.63<math>\pm</math>0.13</b>	98.73 $\pm$ 0.06	<b>51.93<math>\pm</math>0.67</b>	<b>67.42<math>\pm</math>0.52</b>



# Adversarial Robustness



# Qualitative study on generalization capability



# Conclusion

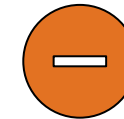
Main Claim:

*“Using Meta-Dropout to perturb the latent features of training examples in a Meta-Learning Framework improves generalization capabilities”*



*Improves:*

- *Decision boundary*
- *Adversarial robustness*
- *Few-Shot learning performance*
- *Hypothesis supported by experiments across large variate of baseline models*
- *Code available*



- *Evaluated on only two datasets*
- *More Shots / More Ways*
- *Relatively small performance increase*
- *Discrepancy in the results*
- *Generalization across datasets domains not discussed*
- *No Computational Cost Reported*
- *Comparison noise  $\phi$  for every layer or shared for all*

# References

- [1]** Hae Beom Lee, Taewook Nam, Eunho Yang, and Sung Ju Hwang. Meta dropout: Learning to perturb latent features for generalization. In ICLR, 2020
- [2]** Alessandro Achille and Stefano Soatto. Information Dropout: Learning Optimal Representations Through Noisy Computation. In PAMI, 2018.
- [3]** Alex Alemi, Ian Fischer, Josh Dillon, and Kevin Murphy. Deep variational information bottleneck. In ICLR, 2017.
- [4]** Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In ICML, 2017.
- [5]** D. P. Kingma, T. Salimans, and M. Welling. Variational Dropout and the Local Reparameterization Trick. In NIPS, 2015
- [6]** Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few shot learning. arXiv preprint arXiv:1707.09835, 2017.
- [7]** Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In NIPS, 2017.
- [8]** Oriol Vinyals, Charles Blundell, Timothy P. Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In NIPS, 2016.