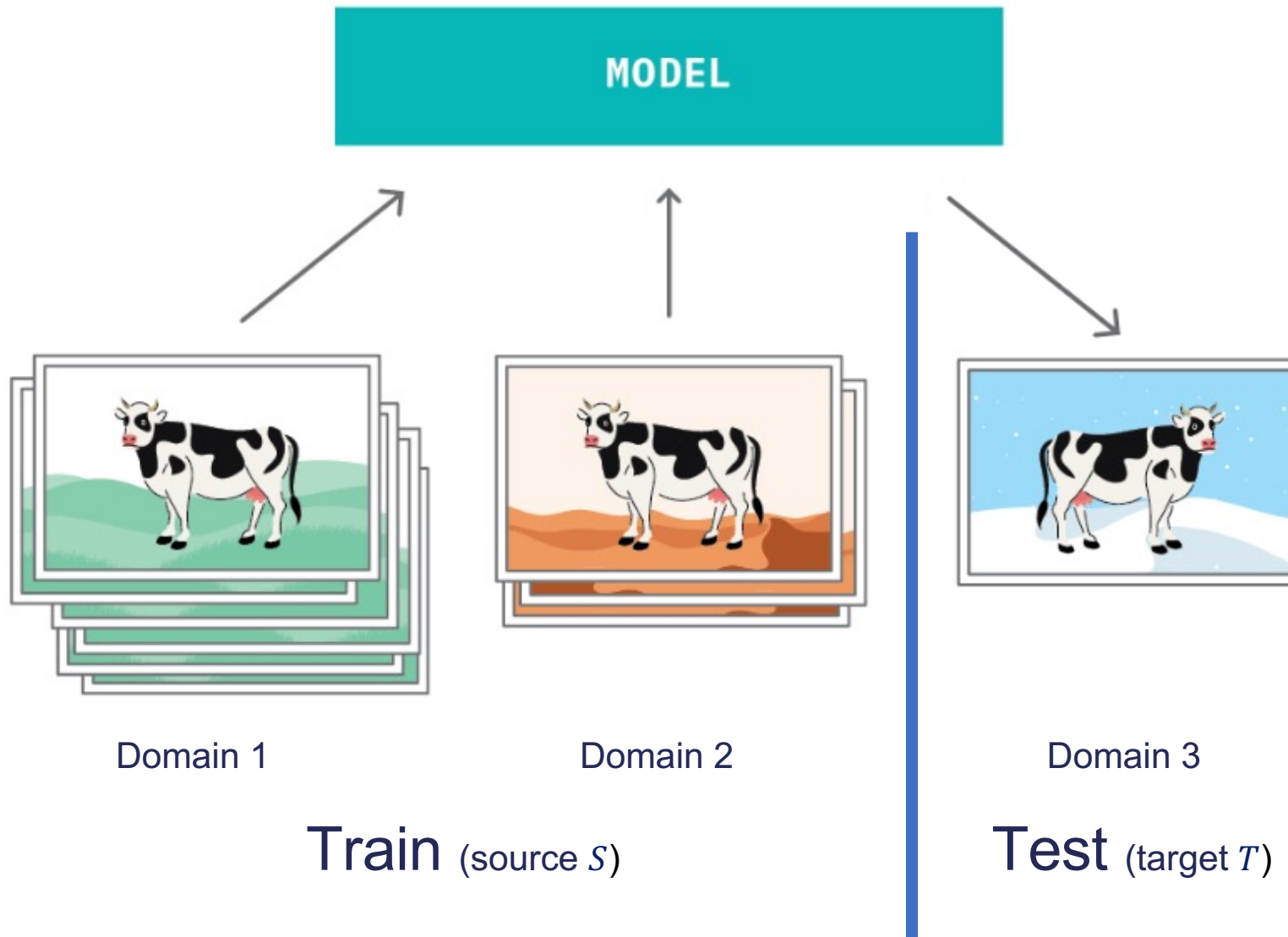


Diverse Weight Averaging for Out-of-Distribution Generalization

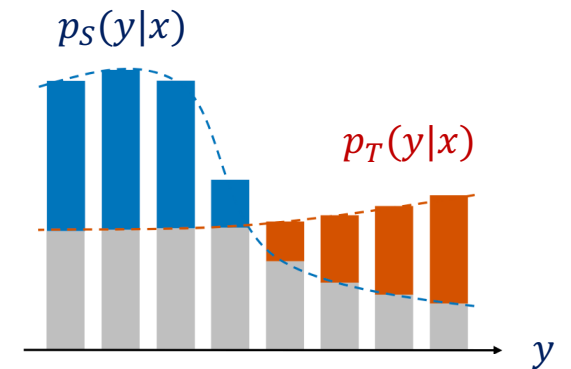
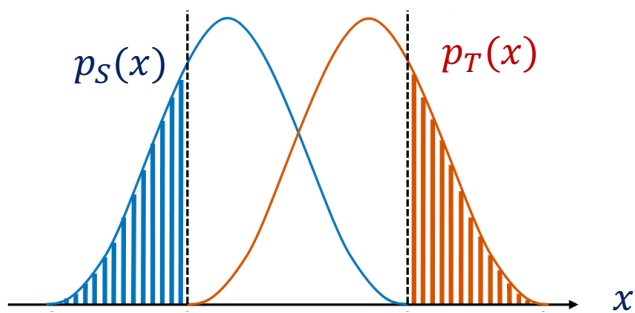
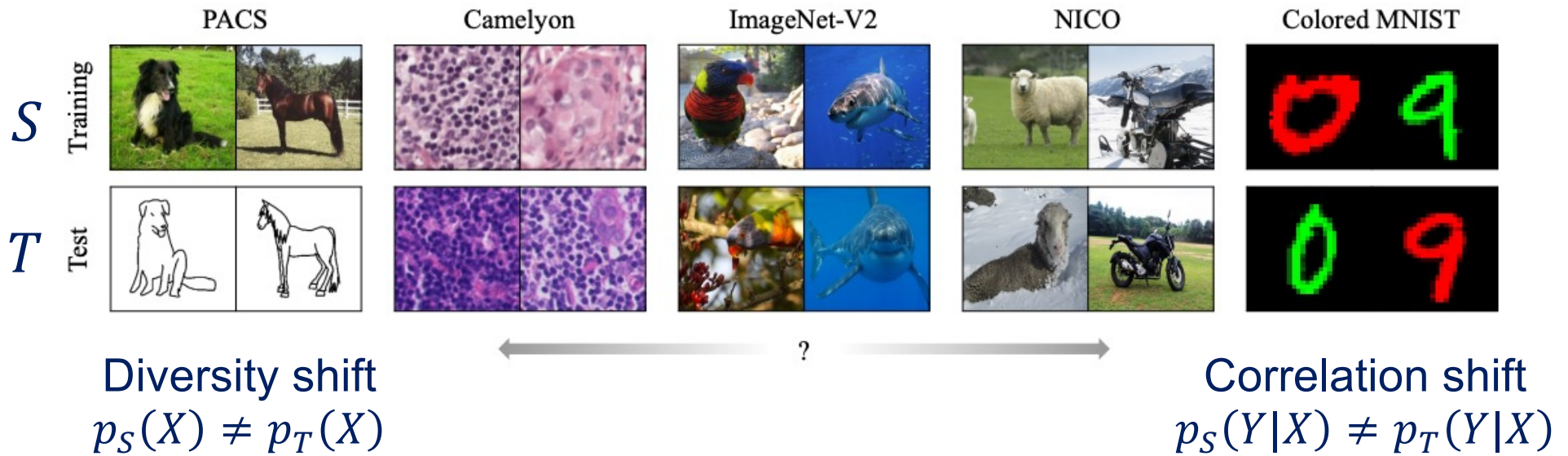
Alexandre Ramé (Sorbonne & Meta AI)
Matthieu Kirchmeyer (Sorbonne & Criteo)
Thibaud Rahier (Criteo)
Alain Rakotomamonjy (LITIS & Criteo)
Patrick Gallinari (Sorbonne & Criteo)
Matthieu Cord (Sorbonne & Valeo.ai)



Goal: generalization to unseen domains



➤ Two kind of source/target distribution shifts





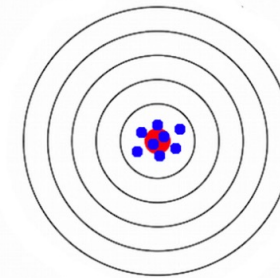
A bias-variance analysis in OOD

$$\mathbb{E}_{\theta} err_T(\theta) = bias^2 + var$$

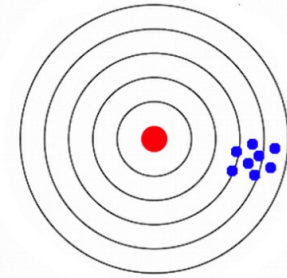
where, with $\bar{f}(x) = \mathbb{E}_{\theta} f_{\theta}(x)$:

- $bias(x, y) = y - \bar{f}(x)$,
- $var(x) = \mathbb{E}_{\theta} \left[\left(f_{\theta}(x) - \bar{f}(x) \right)^2 \right]$.

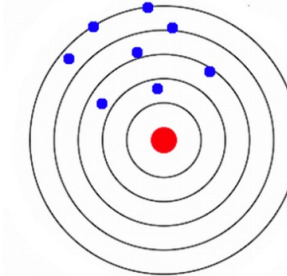
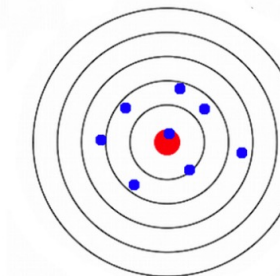
Low
variance



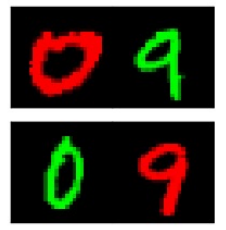
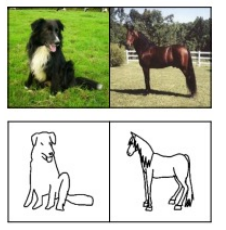
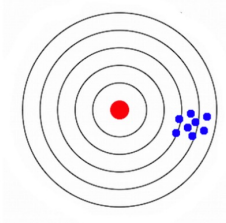
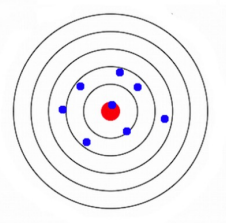
High bias



High
variance

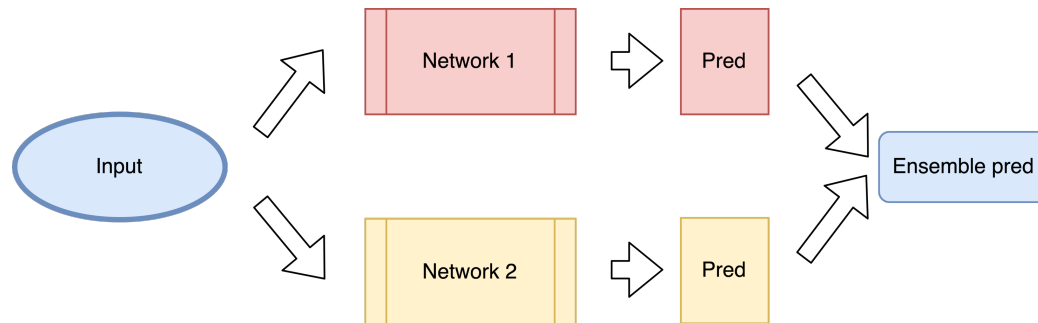


Question: how do *bias* and *var* change with distribution shifts ?

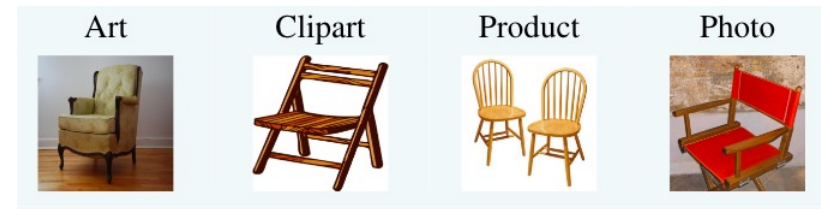
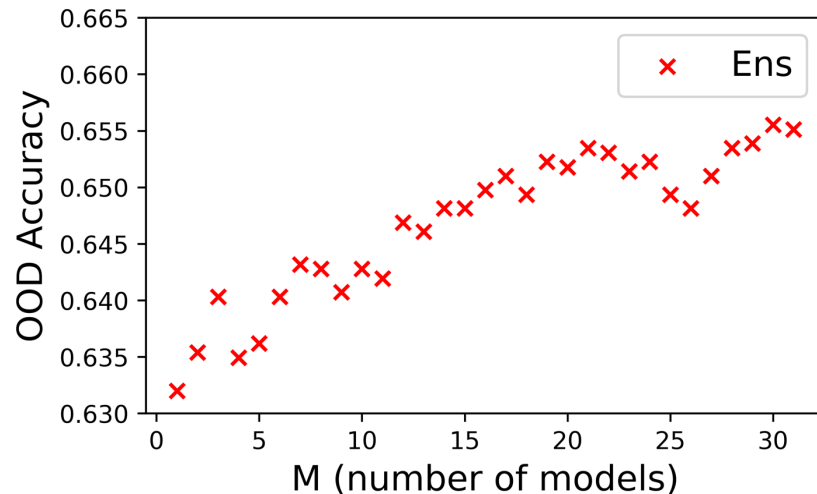
	Correlation shift	Diversity shift
Probabilistic perspective	$p_S(Y X) \neq p_T(Y X)$	$p_S(X) \neq p_T(X)$
Example		
Datasets	ColoredMNIST, CelebA...	OfficeHome, PACS ...
Bias-variance	<p>Large bias Small variance</p> 	<p>Small bias Large variance</p> 
Approaches	<ul style="list-style-type: none"> Invariance: IRM, Coral Robust optimization: gDRO 	<ul style="list-style-type: none"> Variance reduction: ensembling, DiWA



Ensembling M models tackles variance



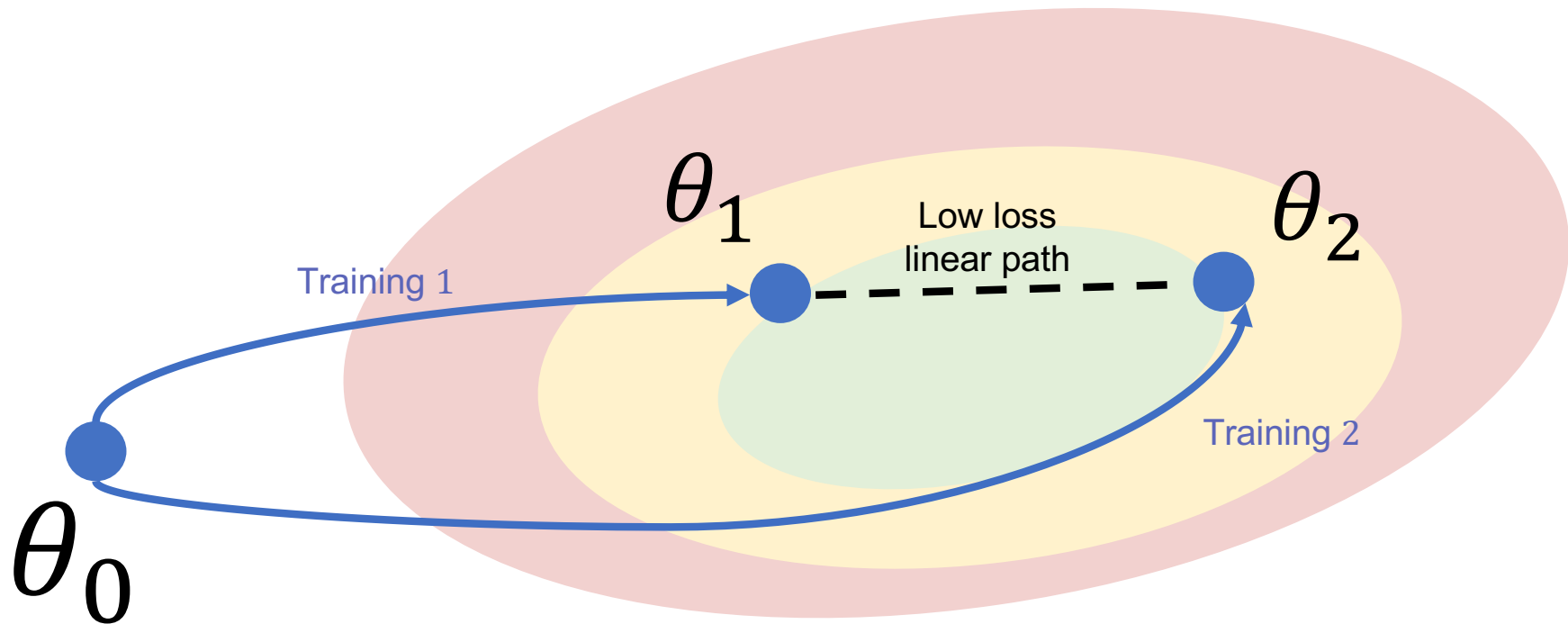
$$\mathbb{E}_{ens_M} err_T(ens_M) \approx bias^2 + \frac{1}{M} var$$



- Setup: OfficeHome under diversity shift
- train on *Clipart*, *Product*, *Photo*
 - test on OOD *Art*

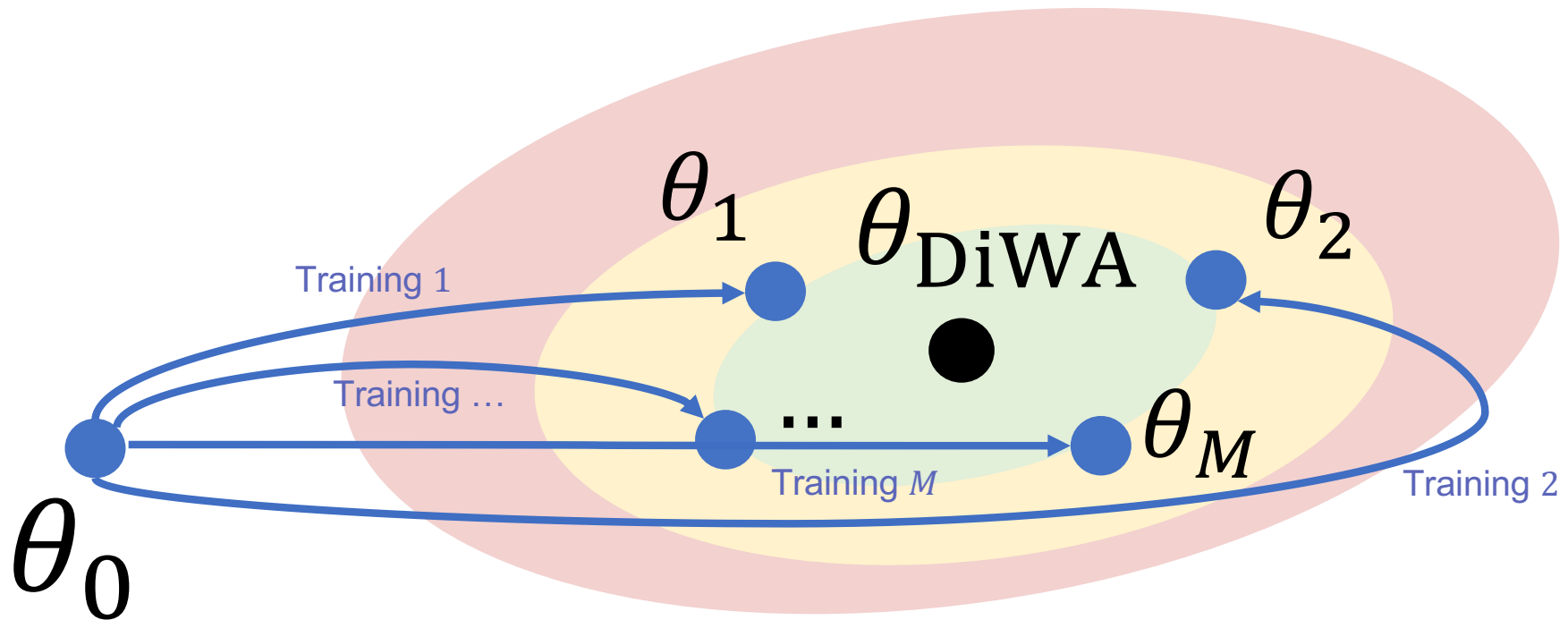
Yet traditional prediction
ensembling is costly ...

➤ An empirical insight: linear mode connectivity



Low-loss linear path
when fine-tunings start from a shared pretrained initialization θ_0
(despite the architecture's non-linearities).

➤ Diverse Weight Averaging (DiWA)



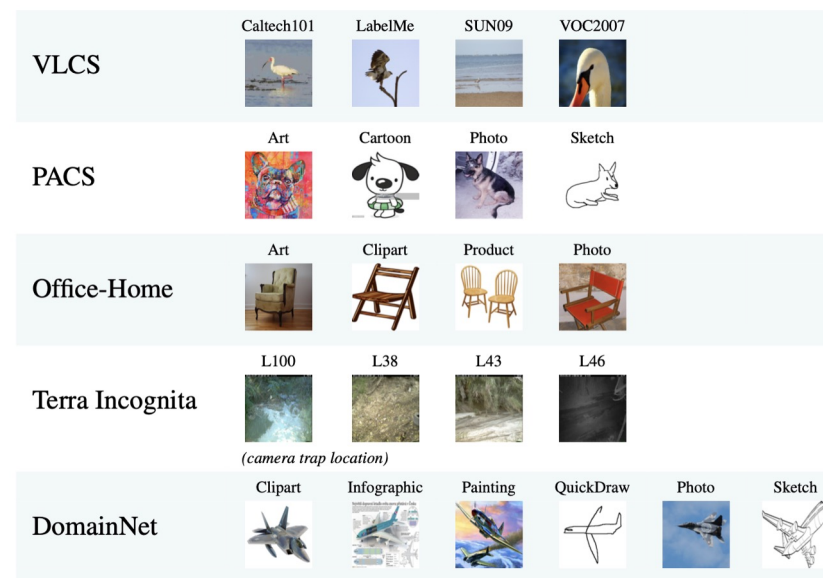
$$\theta_{\text{DiWA}} = \frac{1}{M} \sum_{m=1}^M \theta_m$$

obtained from a shared pretrained initialization θ_0 . Then:

$$f_{\theta_{\text{DiWA}}} = f_{\frac{1}{M} \sum_{m=1}^M \theta_m} \approx \frac{1}{M} \sum_{m=1}^M f_{\theta_m}.$$

➤ SoTA on DomainBed [Gulrajani2021]

Reference benchmark for OOD generalization in computer vision, imposing the *code, datasets, training procedures, hyperparameter search* etc.



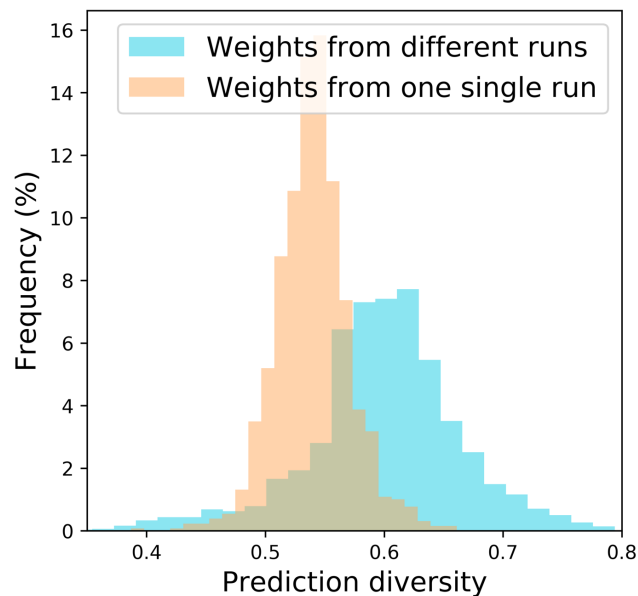
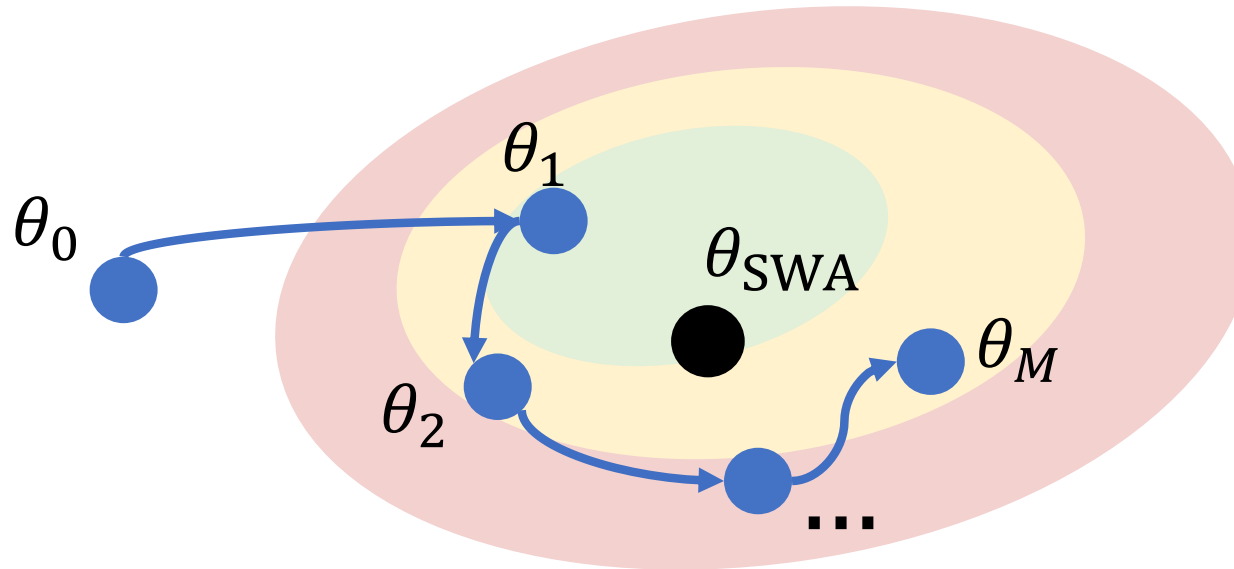
Algo	Cost	PACS	VLCS	OH	TI	DN	Avg
ERM	1	85.5	77.5	66.5	46.1	40.9	63.3
CORAL	1	86.2	78.8	68.7	47.6	41.5	64.6
SWAD	1	88.1	79.1	70.6	50.0	46.5	66.9
ENS	20	88.1	78.5	71.7	50.8	47.0	67.2
DiWA	1	89.0	78.6	72.8	51.9	47.7	68.0

[Gulrajani2021] In search of lost domain generalization. ICLR

[Cha2021] SWAD: Domain Generalization by Seeking Flat Minima. NeurIPS



Previous SoTA: single-run WA



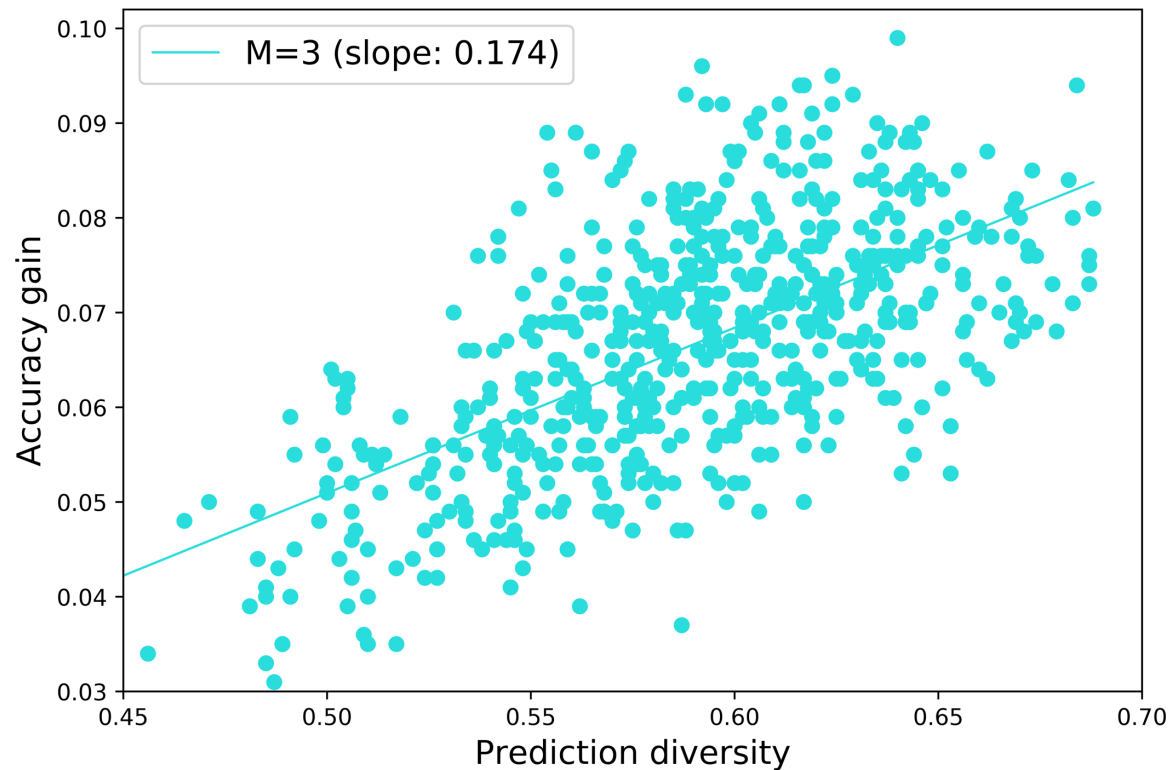
Weights from different runs are more diverse (left) thus their average is better (next slide).



Covariance as diversity

$$\mathbb{E}_{\theta_{WA}} \text{err}_T(\theta_{WA}) \approx \text{bias}^2 + \frac{1}{M} \text{var} + \frac{M-1}{M} \text{cov},$$

where *cov* is smaller when models are uncorrelated, *i.e.*, functionally diverse.

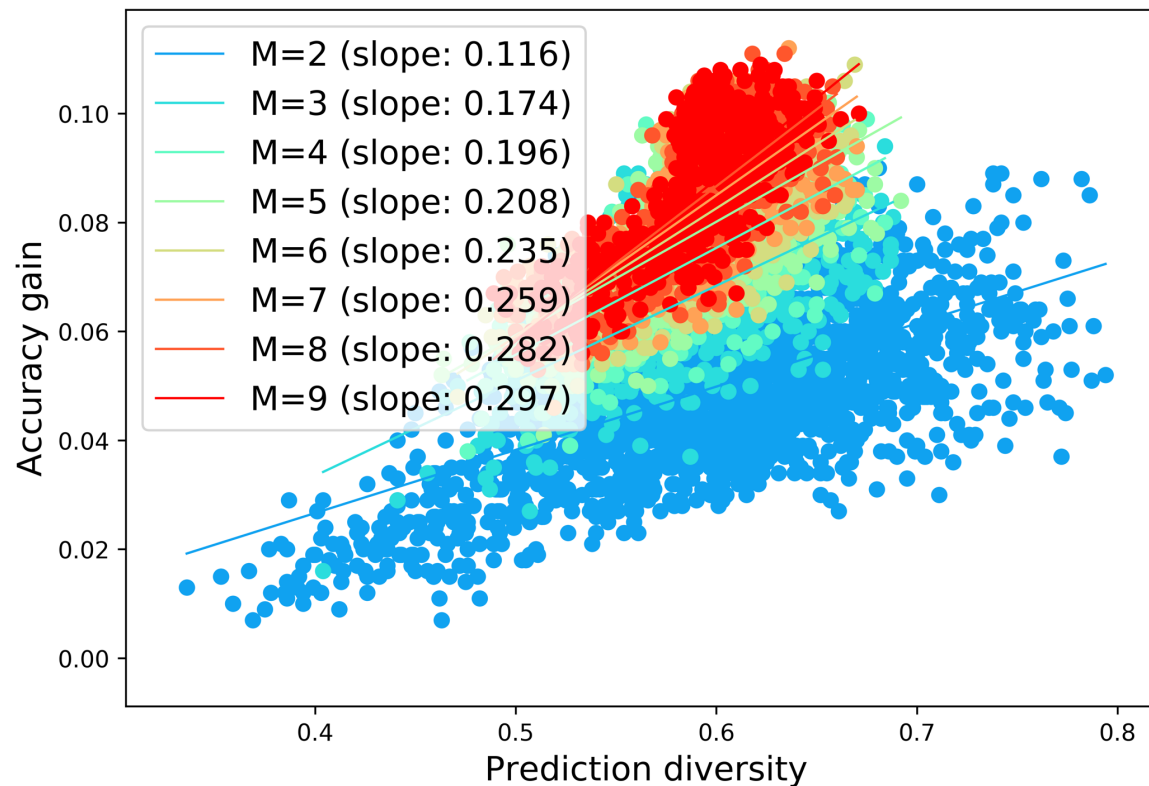


Legend: Each dot is the accuracy gain of averaging $M = 3$ models over the average accuracy *wrt* their diversity (normalized count of different errors).



Diversity even more important for larger M

$$\mathbb{E}_{\theta_{WA}} \text{err}_T(\theta_{WA}) = \text{bias}^2 + \frac{1}{M} \text{var} + \frac{M-1}{M} \text{cov} + \mathcal{O}(\bar{\Delta}^2).$$



Legend: Each dot is the accuracy gain of averaging M models over the average accuracy wrt their diversity (normalized count of different errors).



Conclusion

- ❖ Bias-variance analysis in OOD
 - ✓ Relate diversity shift to variance
 - ✓ Relate correlation shift to bias
- ❖ New weight averaging strategy
 - ✓ Average all weights obtained from the hyperparameter search
 - ✓ SoTA on DomainBed to tackle diversity shift

arXiv: <https://arxiv.org/abs/2205.09739>

code: <https://github.com/alexrame/diwa>