

Real-World Learning

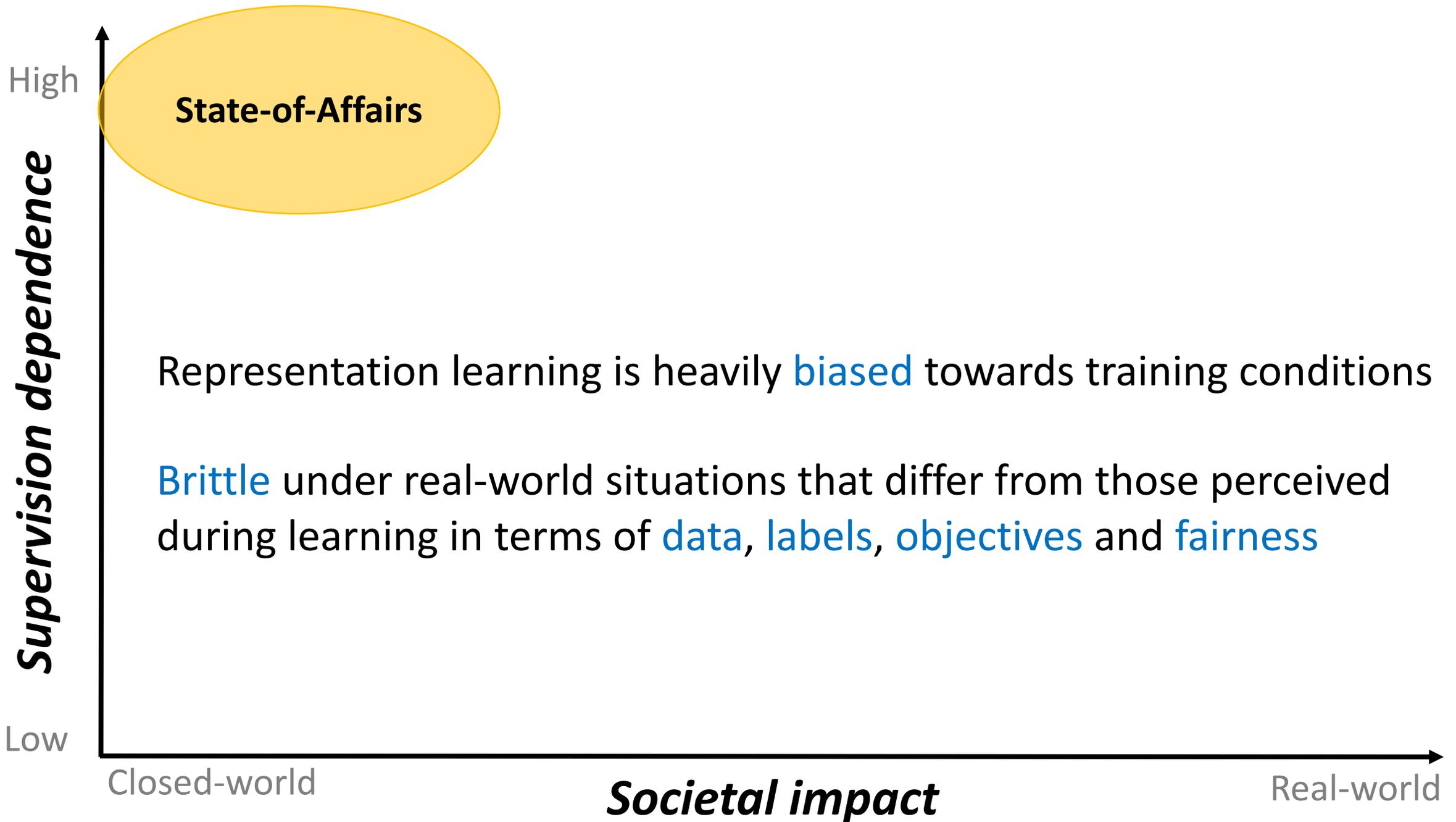
Prof. dr. Cees Snoek

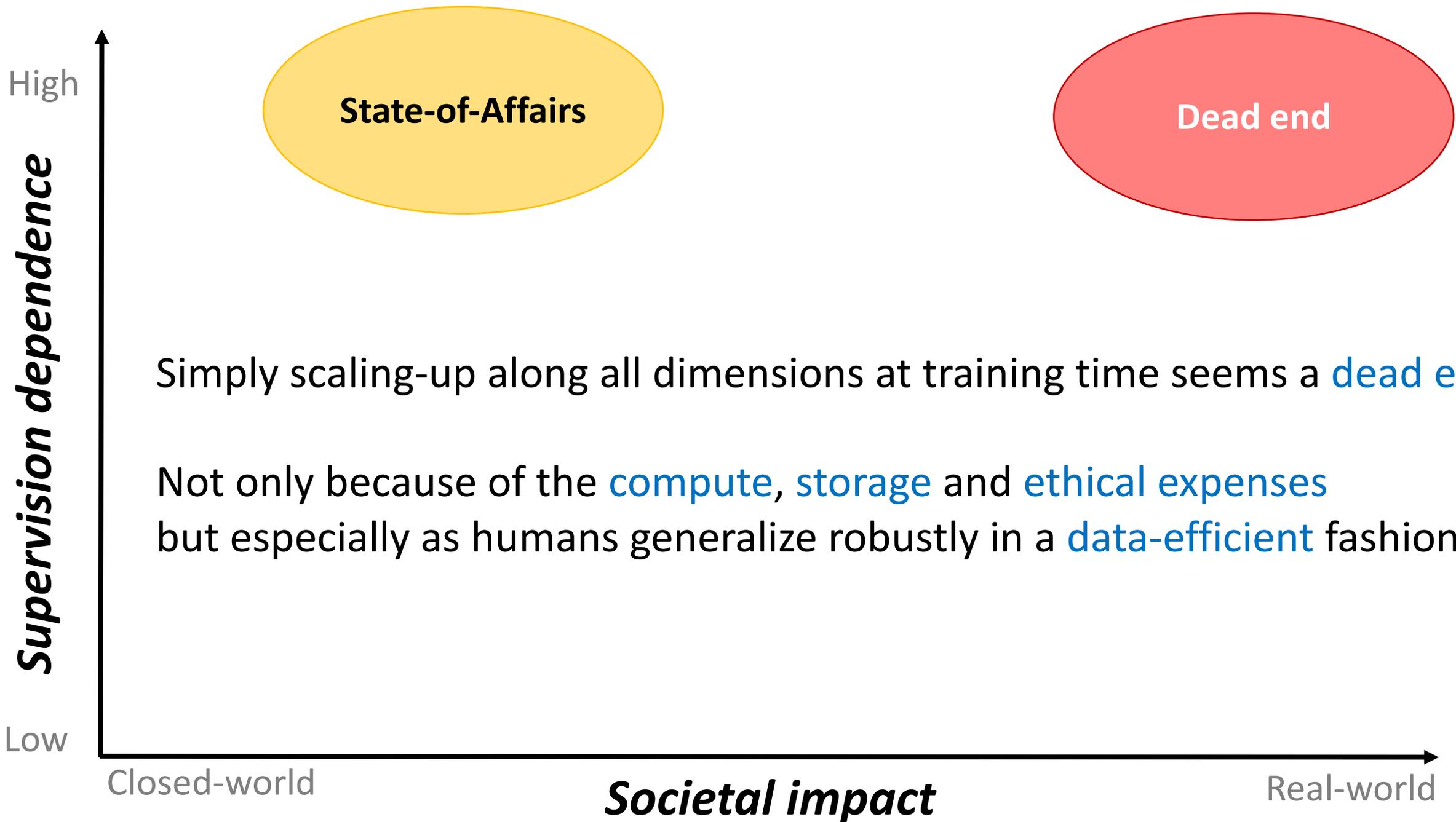
Video & Image Sense Lab



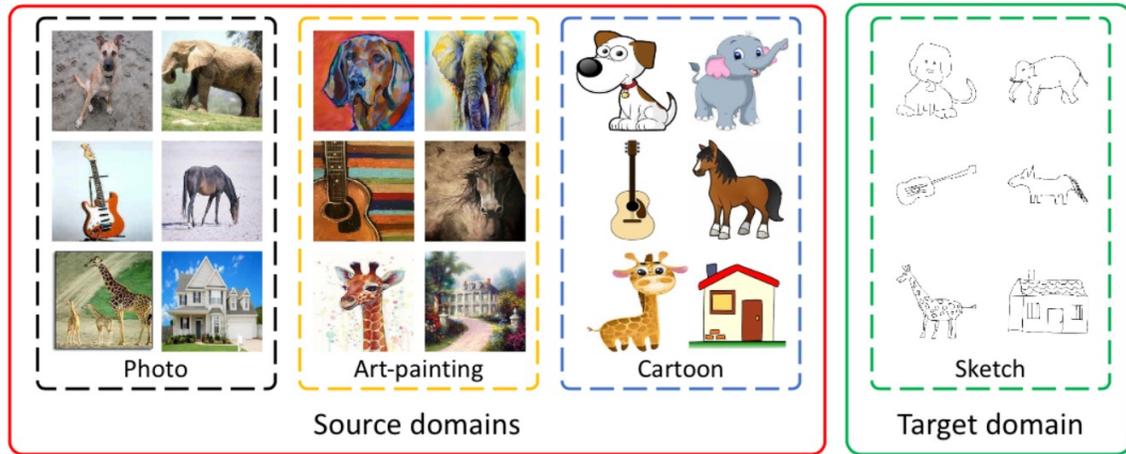
Awesome learning





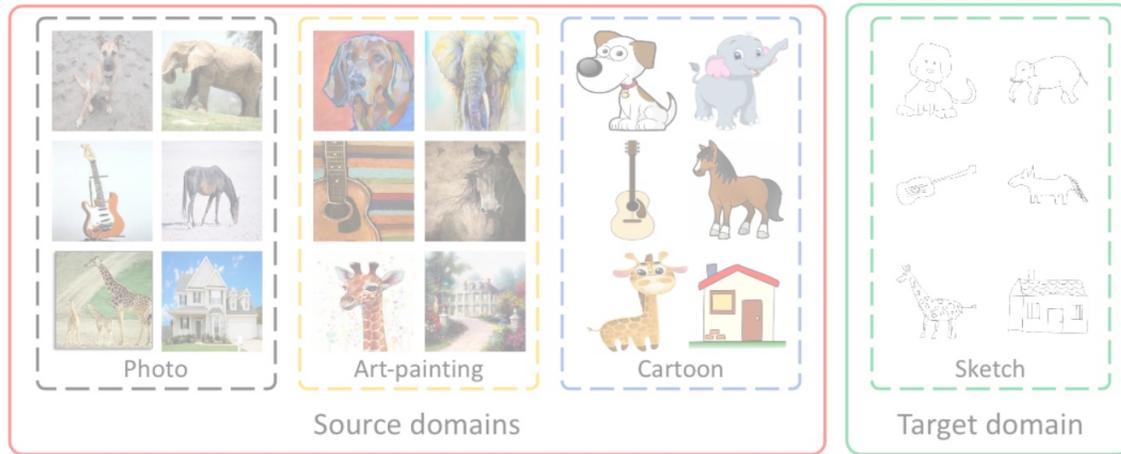


Distribution gap



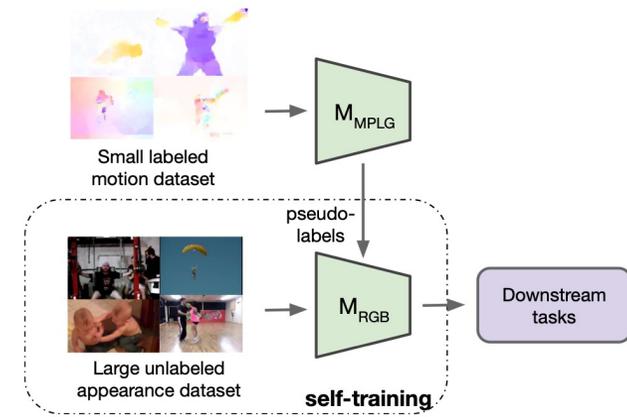
w/ Zehao Xiao *et al.*, ICML 2021

Distribution gap



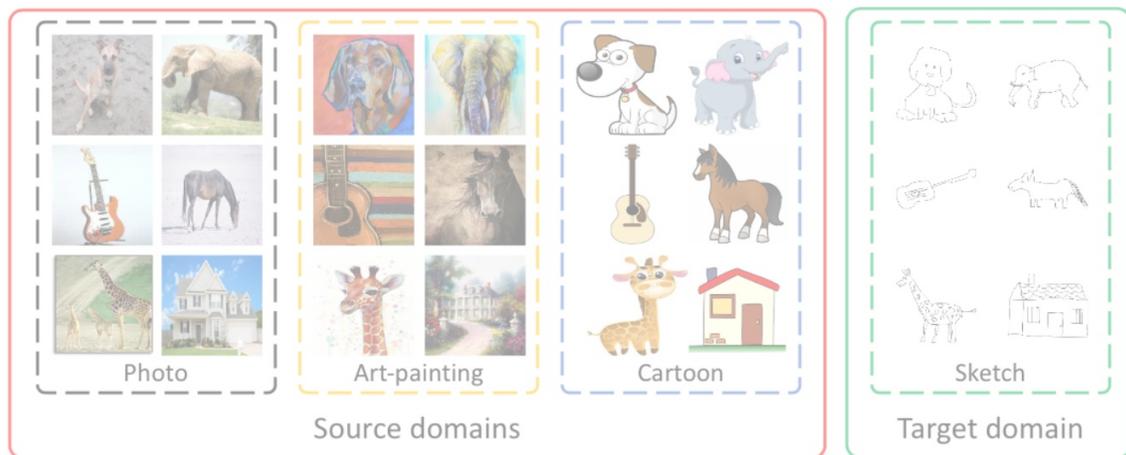
w/ Zehao Xiao *et al.*, ICML 2021

Label gap



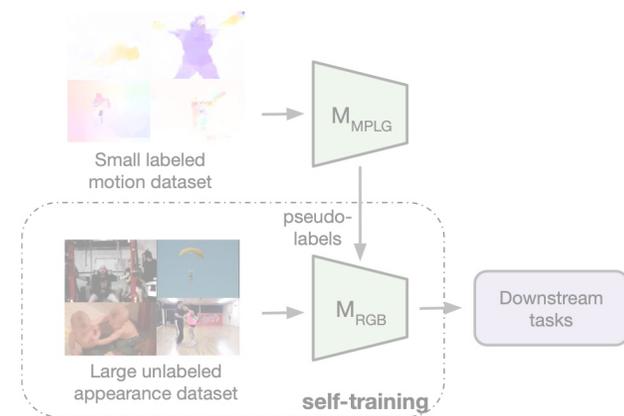
w/ Kirill Gavriluk *et al.*, ICCV 2021

Distribution gap



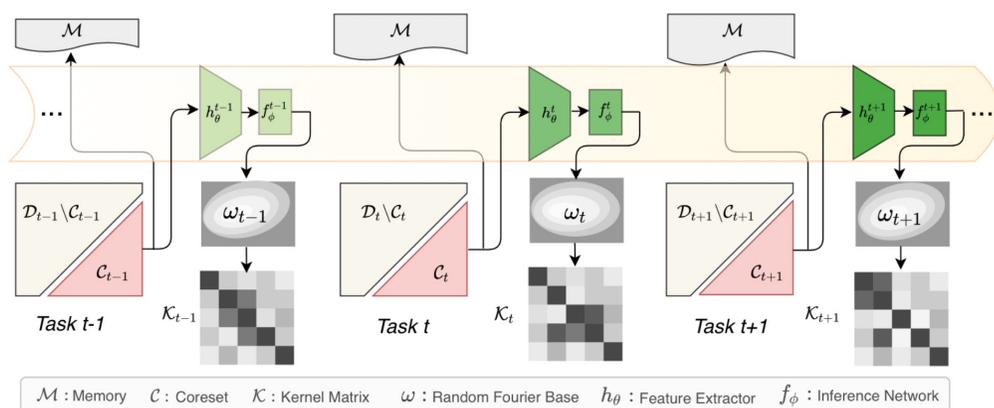
w/ Zehao Xiao *et al.*, ICML 2021

Label gap



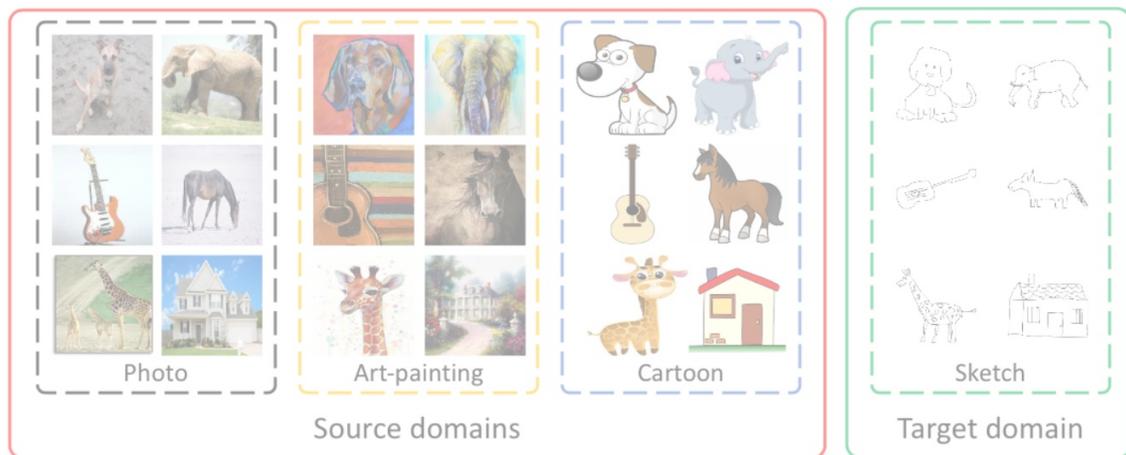
w/ Kirill Gavriluk *et al.*, ICCV 2021

Objectives gap



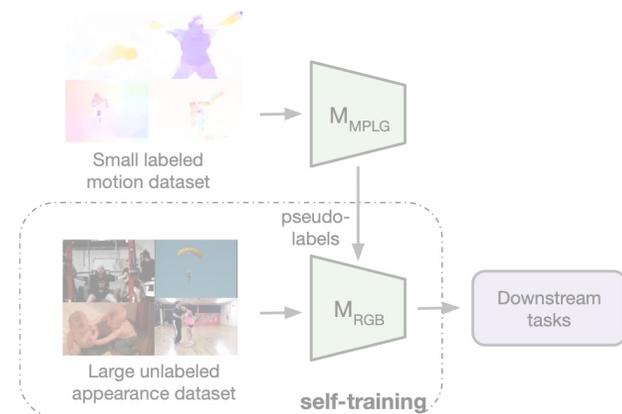
w/ Mohammad Mahdi Derakhshani *et al.*, ICML 2021

Distribution gap



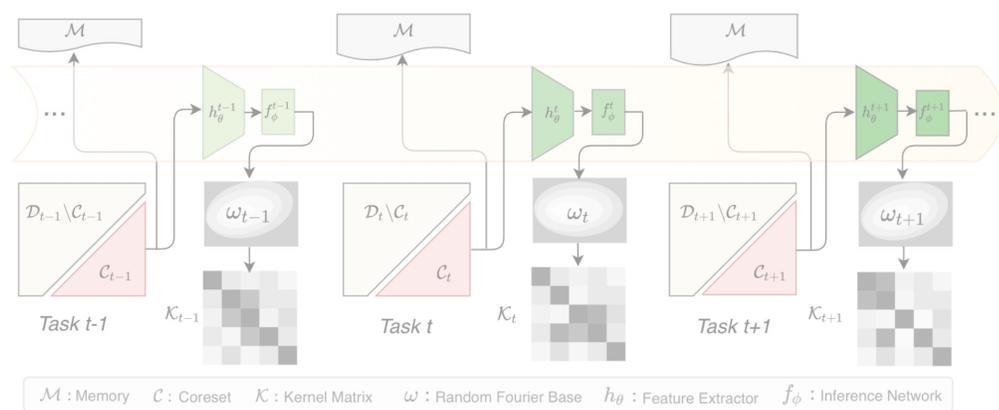
w/ Zehao Xiao *et al.*, ICML 2021

Label gap



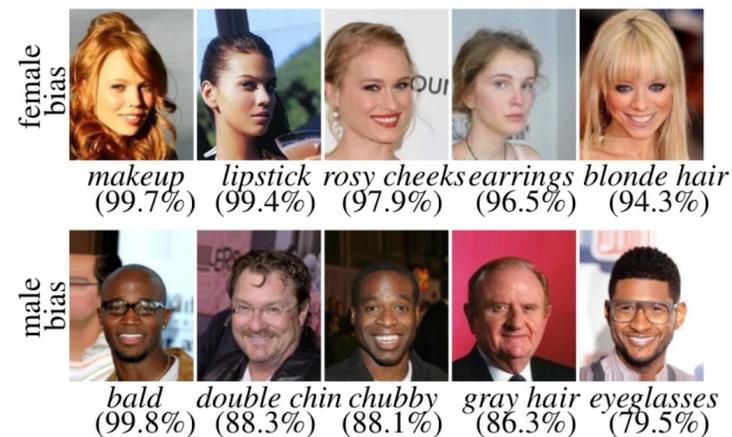
w/ Kirill Gavrilyuk *et al.*, ICCV 2021

Objectives gap

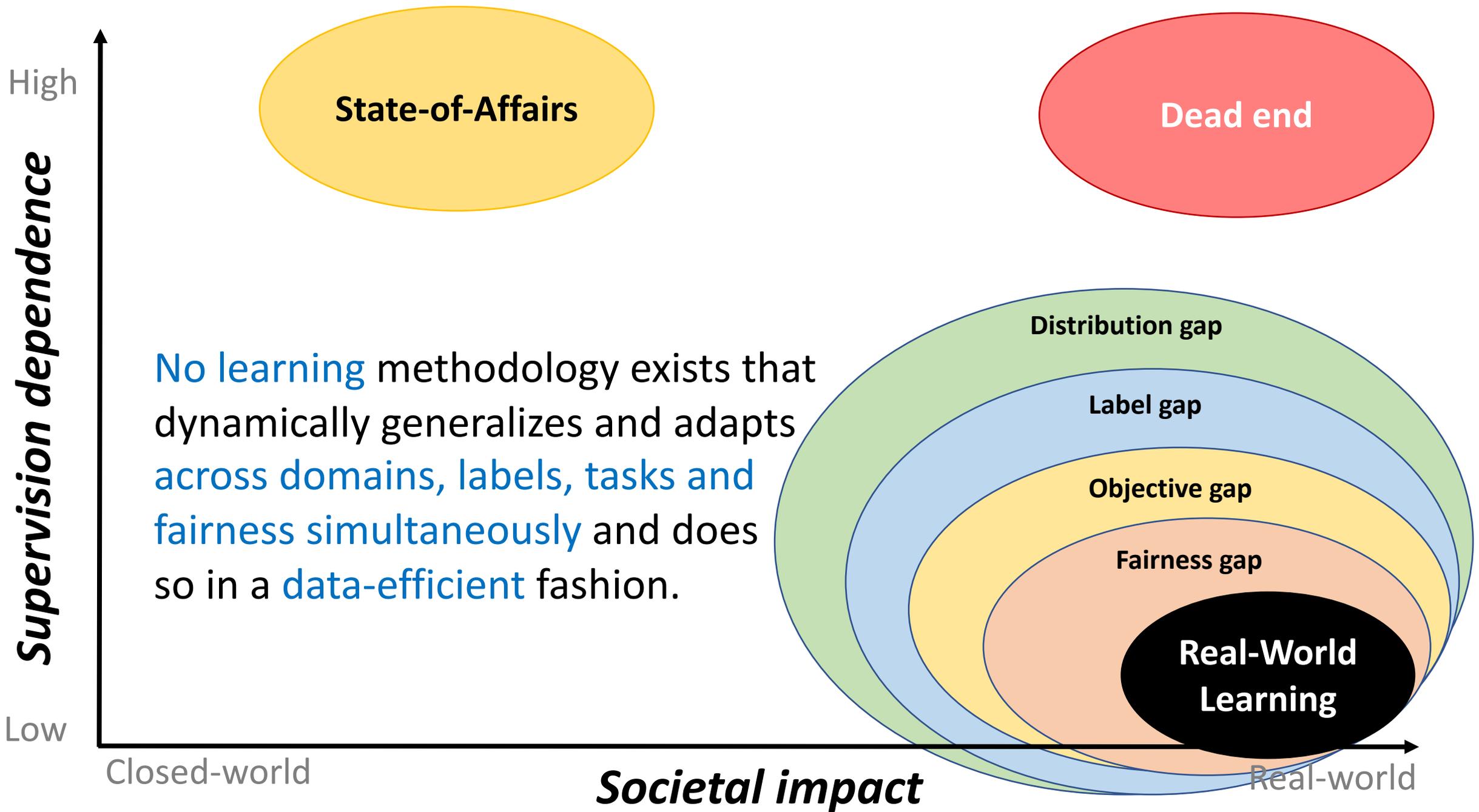


w/ Mohammad Mahdi Derakhshani *et al.*, ICML 2021

Fairness gap



w/ William Thong, BMVC 2021



This talk

We question common representation learning assumptions

i. Learning without **label** assumption

ii. Learning without **task** assumption

iii. Learning without **domain** assumption

1. Learning without label assumption



Pengwan Yang
University of Amsterdam



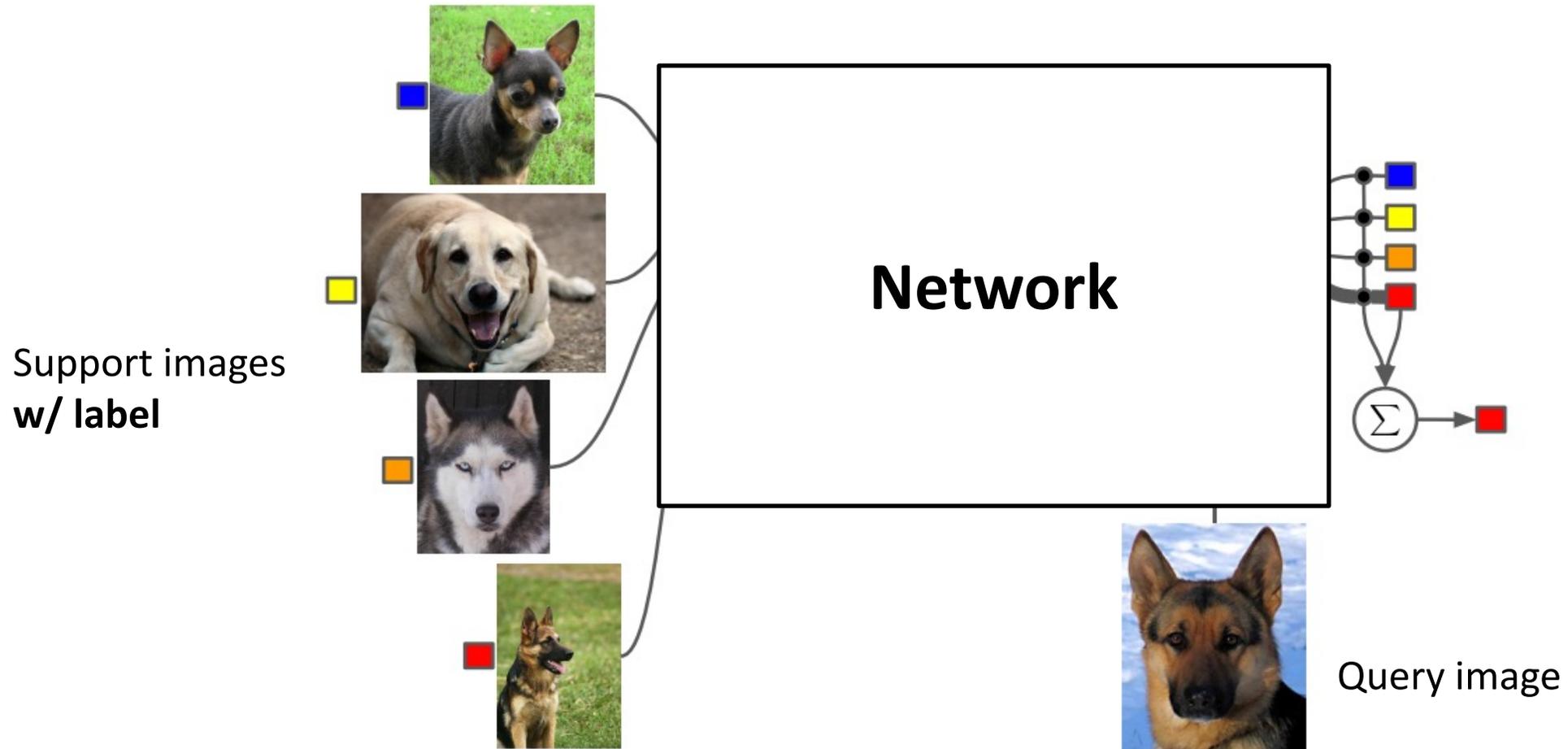
Pascal Mettes
University of Amsterdam



Cees Snoek
University of Amsterdam

Few-Shot Transformation of Common Actions into Time and Space. In *CVPR* 2021.

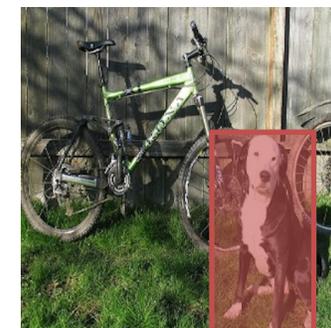
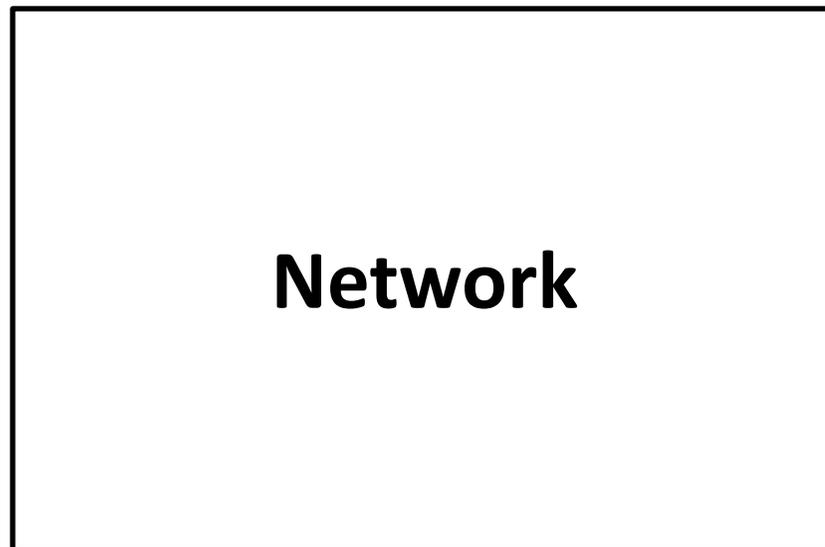
Canonical Paradigm: few-shot classification



Canonical Paradigm: few-shot detection



Support images
w/ label + box

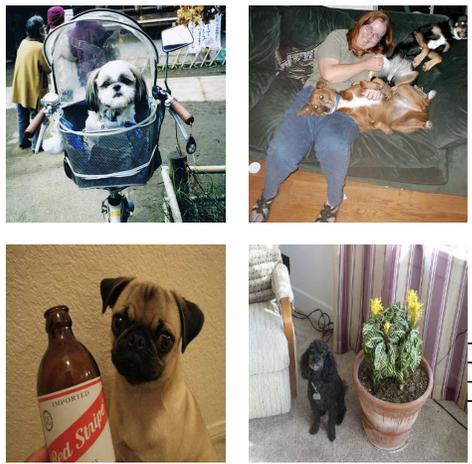


Query image

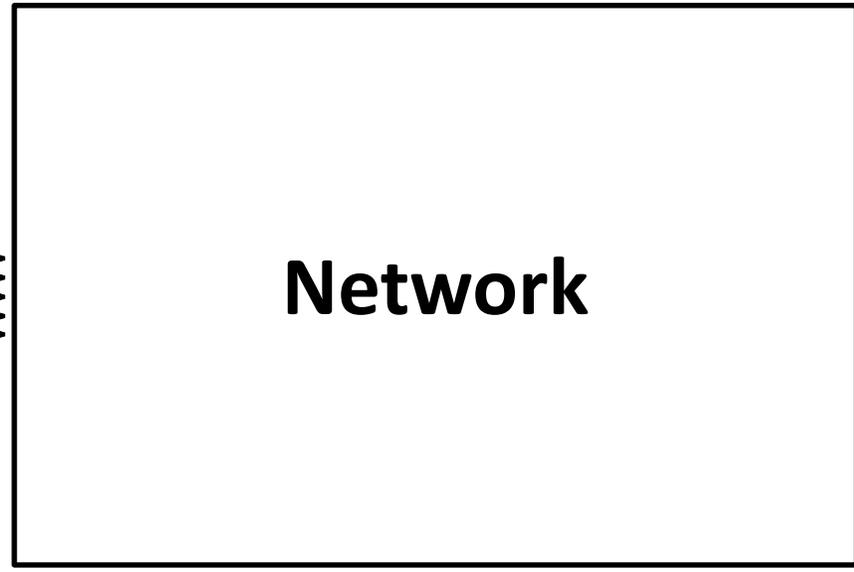


Few-shot common object localization

Support images **w/o label and w/o box**

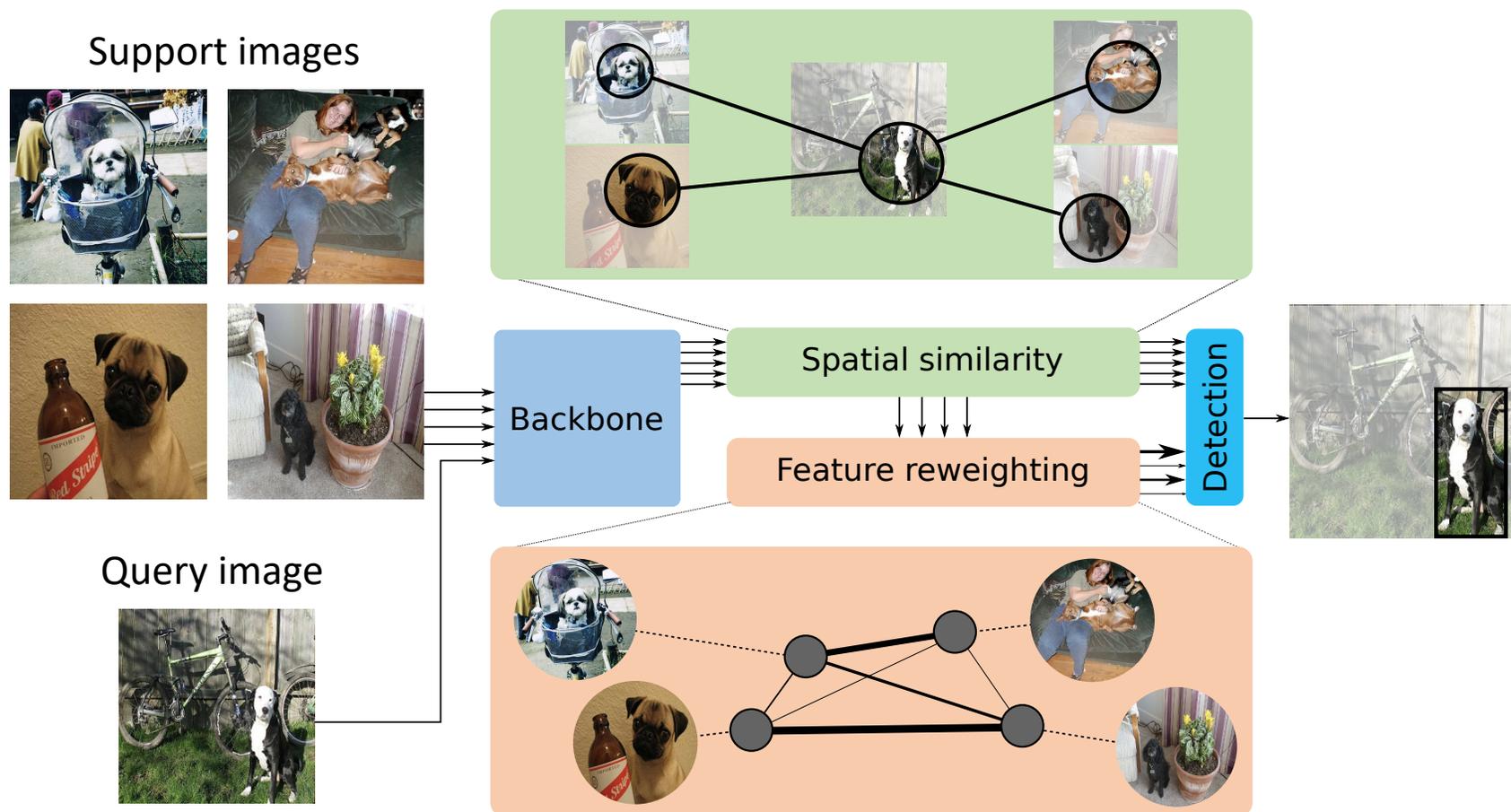


Query image



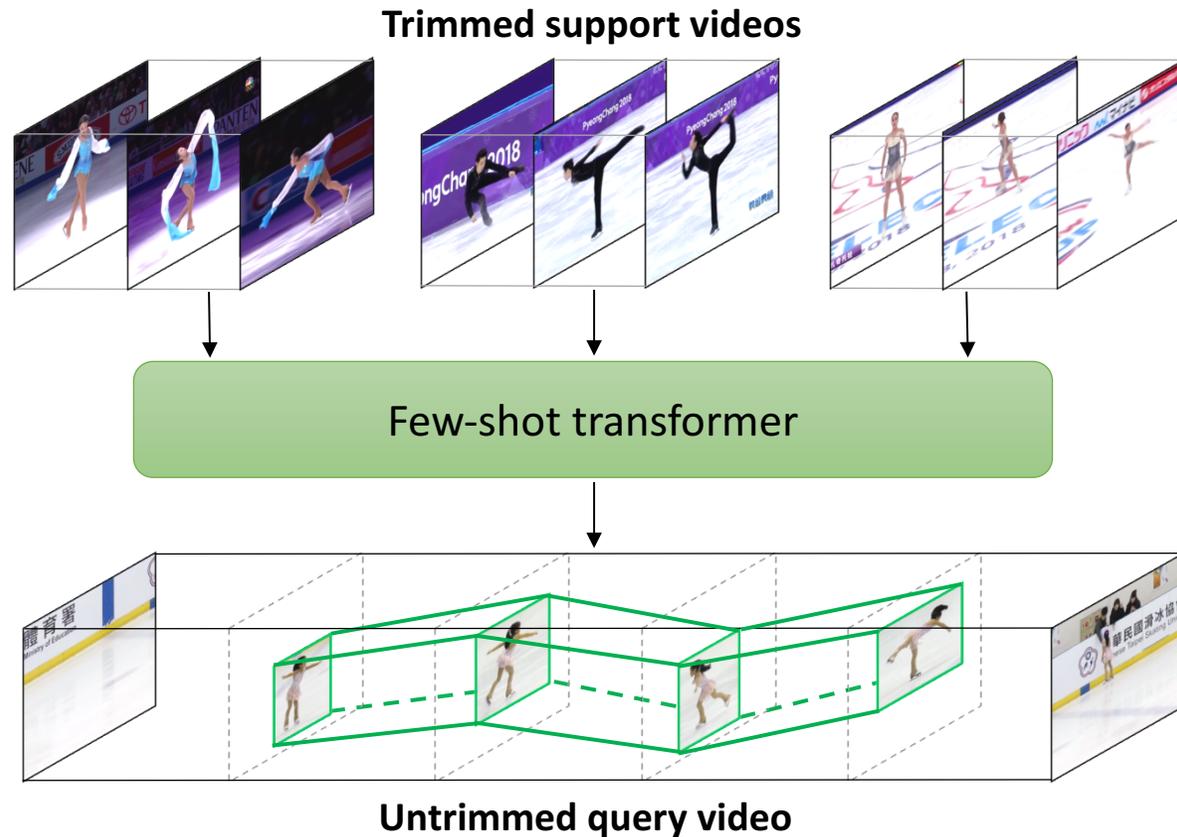
*Localize the common object in the query image **without** any label and box annotation*

Few-shot common object localization



Localize the common object in the query image without any box annotation

Few-shot common action in video



No need for action class label or any temporal and/or spatial annotation

Example

support videos



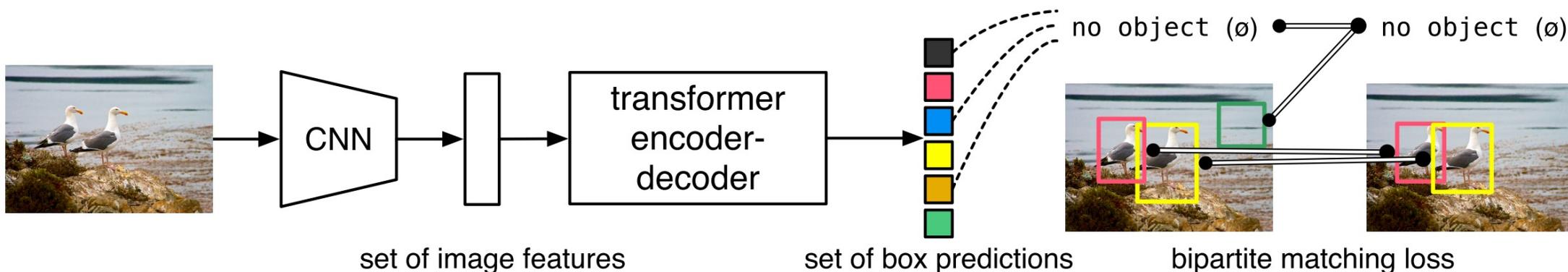
query video



one-shot prediction

five-shot prediction

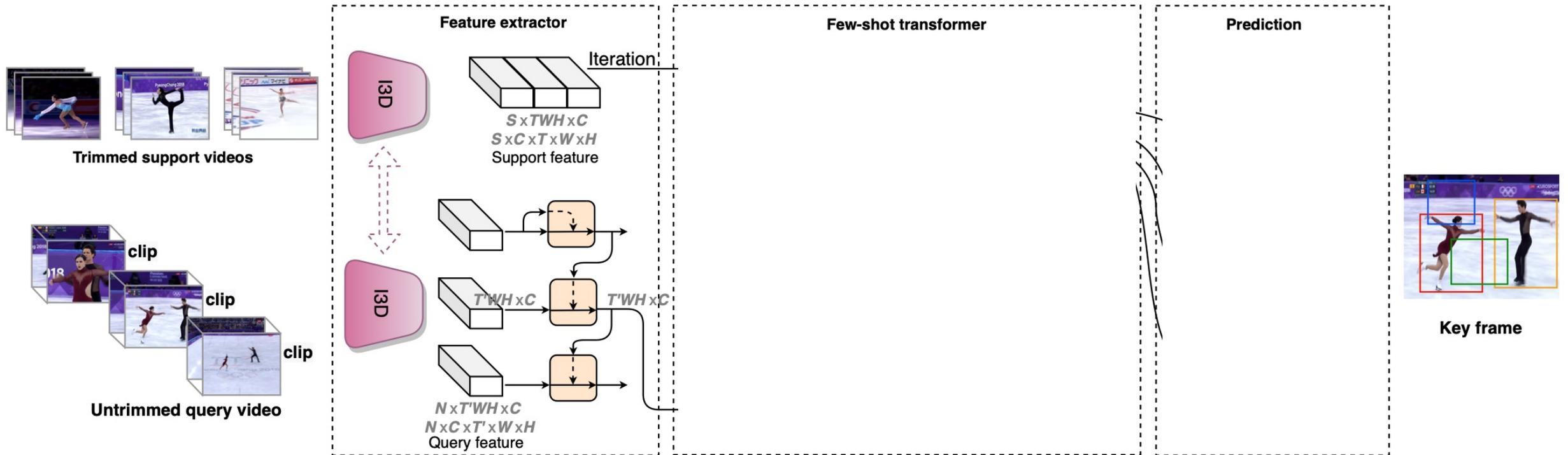
Inspiration: object detection transformers



Benefits of transformers:

- i) it avoids the needle-in-the-haystack problem with proposals
- ii) it provides powerful relation modeling capability

Method



S : number of support videos

T : temporal length of support video feature

H : feature height

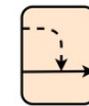
W : feature width

N : number of query clips

T' : temporal length of query clip feature

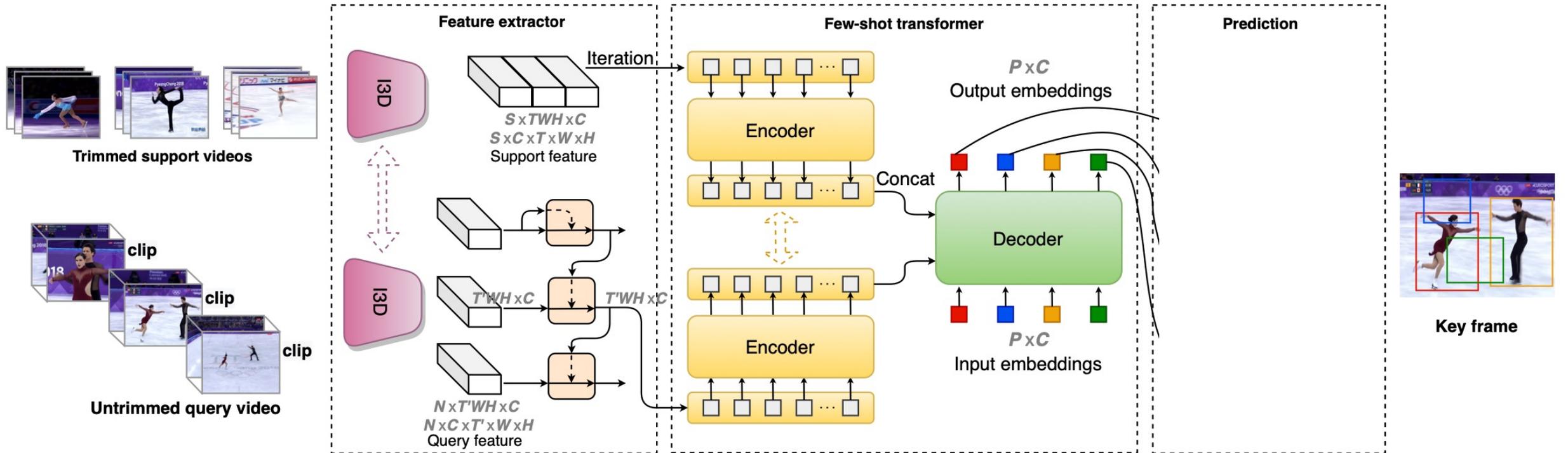
C : feature dimension

P : number of predictions



Common attention block

Method



S : number of support videos

T : temporal length of support video feature

H : feature height

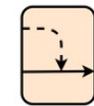
W : feature width

N : number of query clips

T' : temporal length of query clip feature

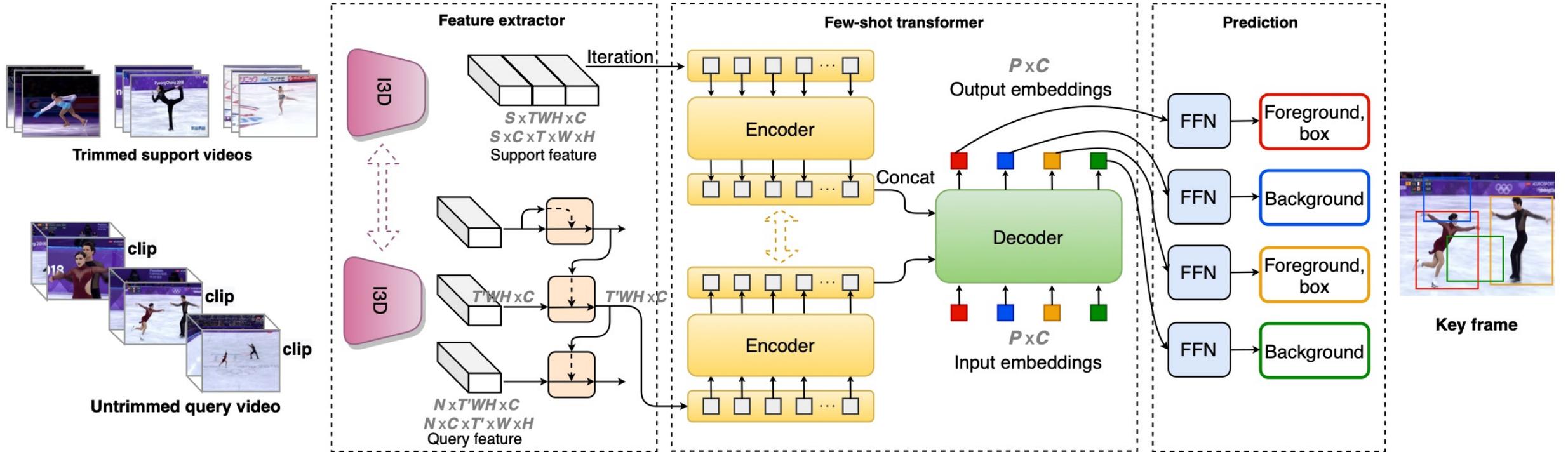
C : feature dimension

P : number of predictions



Common attention block

Method



S : number of support videos

T : temporal length of support video feature

H : feature height

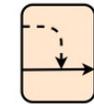
W : feature width

N : number of query clips

T' : temporal length of query clip feature

C : feature dimension

P : number of predictions



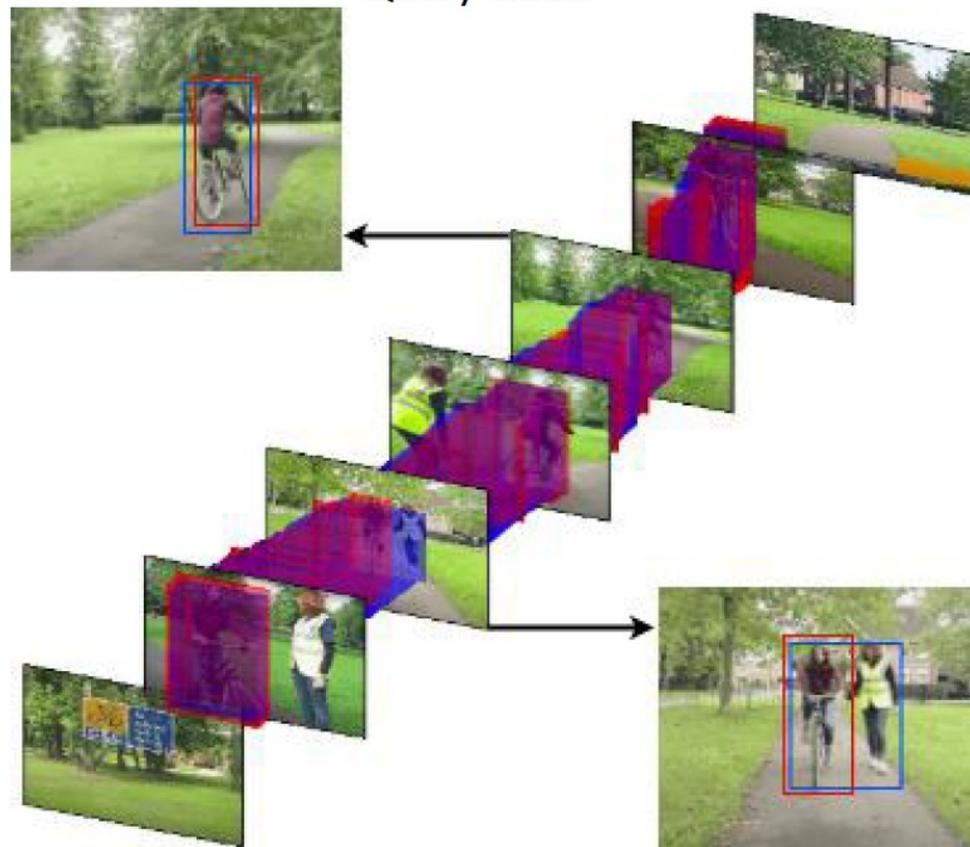
Common attention block

Results

Support videos



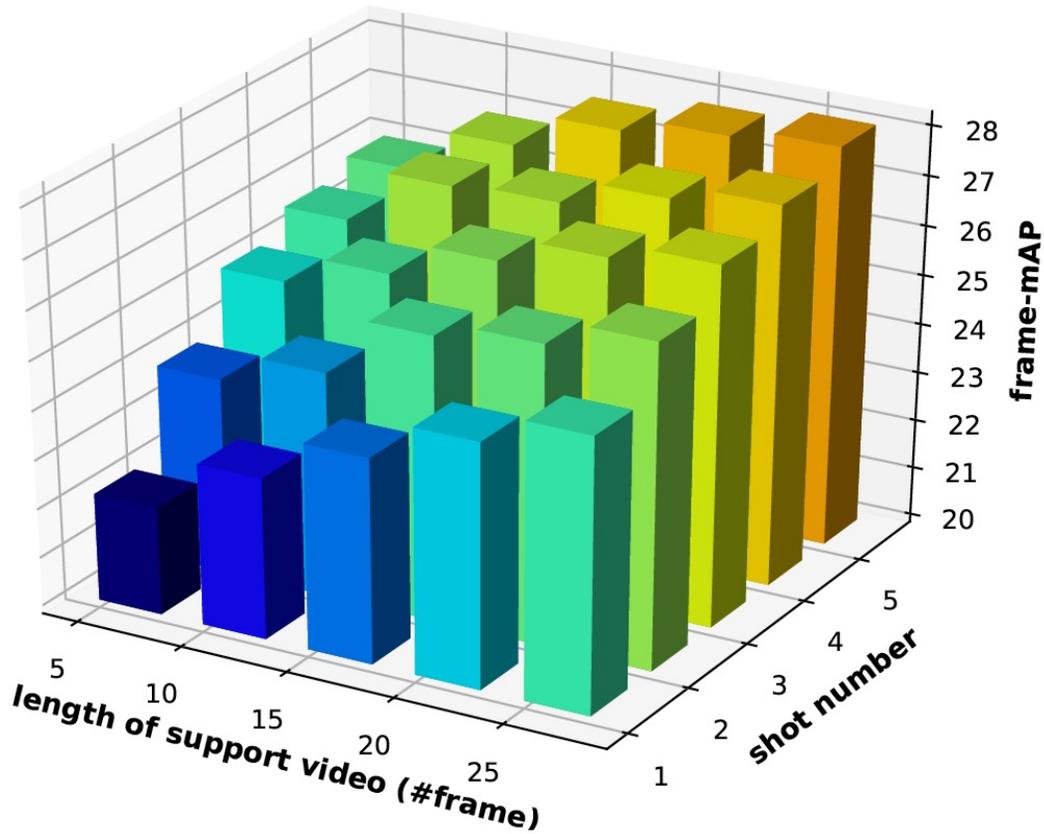
Query video



one-shot (blue)
five-shot (red)

Common action localization in time and space

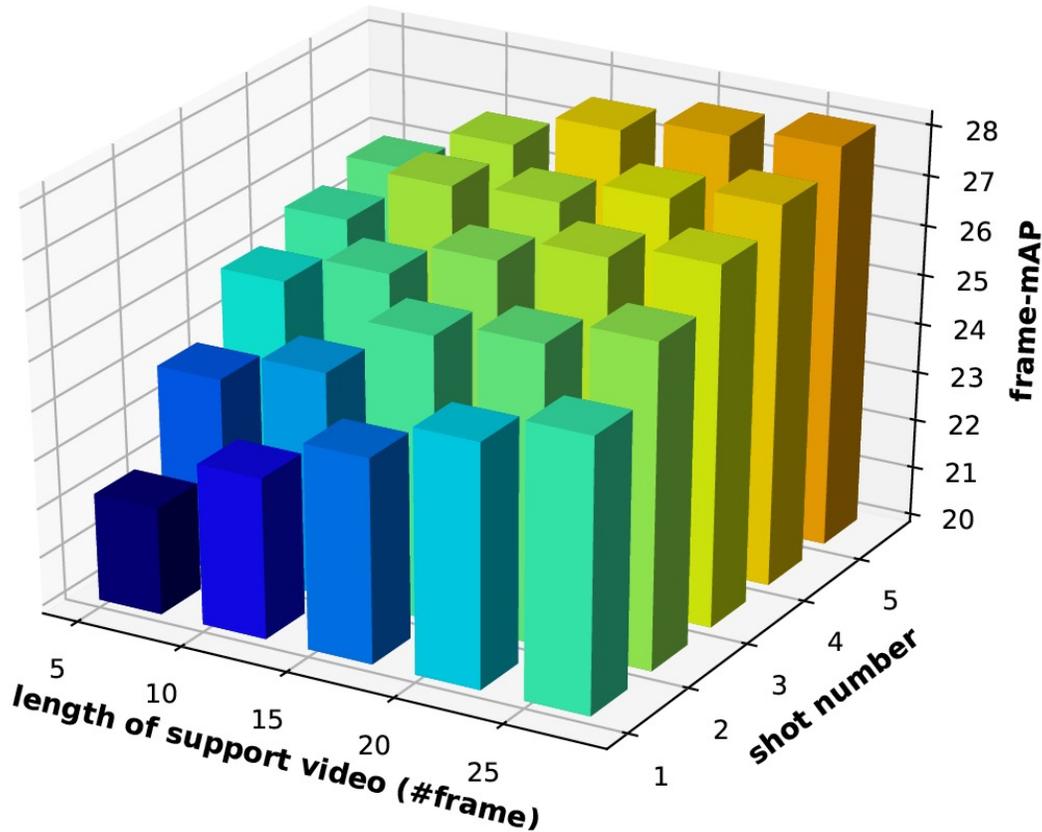
Ablations on Common-AVA



Influence of length and number of support videos.

We obtain a more precise common localization with more and longer support videos.

Ablations on Common-AVA

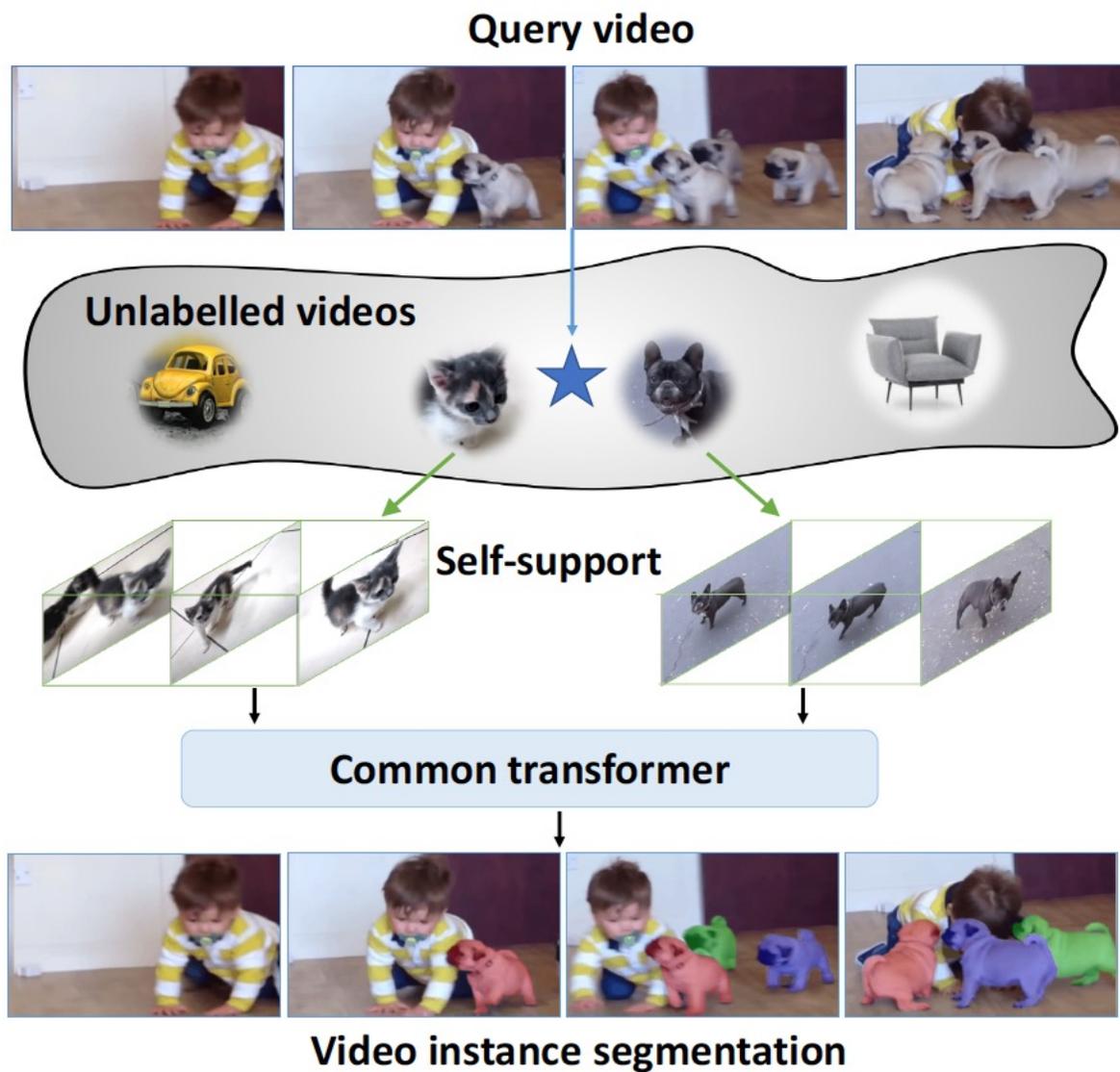


Influence of length and number of support videos.
We obtain a more precise common localization with more and longer support videos.

No noise	28.1
Video-level noise	
1 noisy support video of other class	26.8
1 noisy support video without action	26.3
2 noisy support videos of different class	25.3
2 noisy support videos of same class	24.7
Frame-level noise	
2 noisy frames in each support video	27.9
4 noisy frames in each support video	27.4
6 noisy frames in each support video	26.1
8 noisy frames in each support video	24.5

Effect of noisy support videos for the five-shot setting. The result shows our robustness.

Self-support video instance segmentation



Find support videos using the query

Pool from unlabelled video in self-supervised fashion

Transformer enables instance segmentation

No labels, no masks.

Video instance segmentation results

Self-support examples



Query video

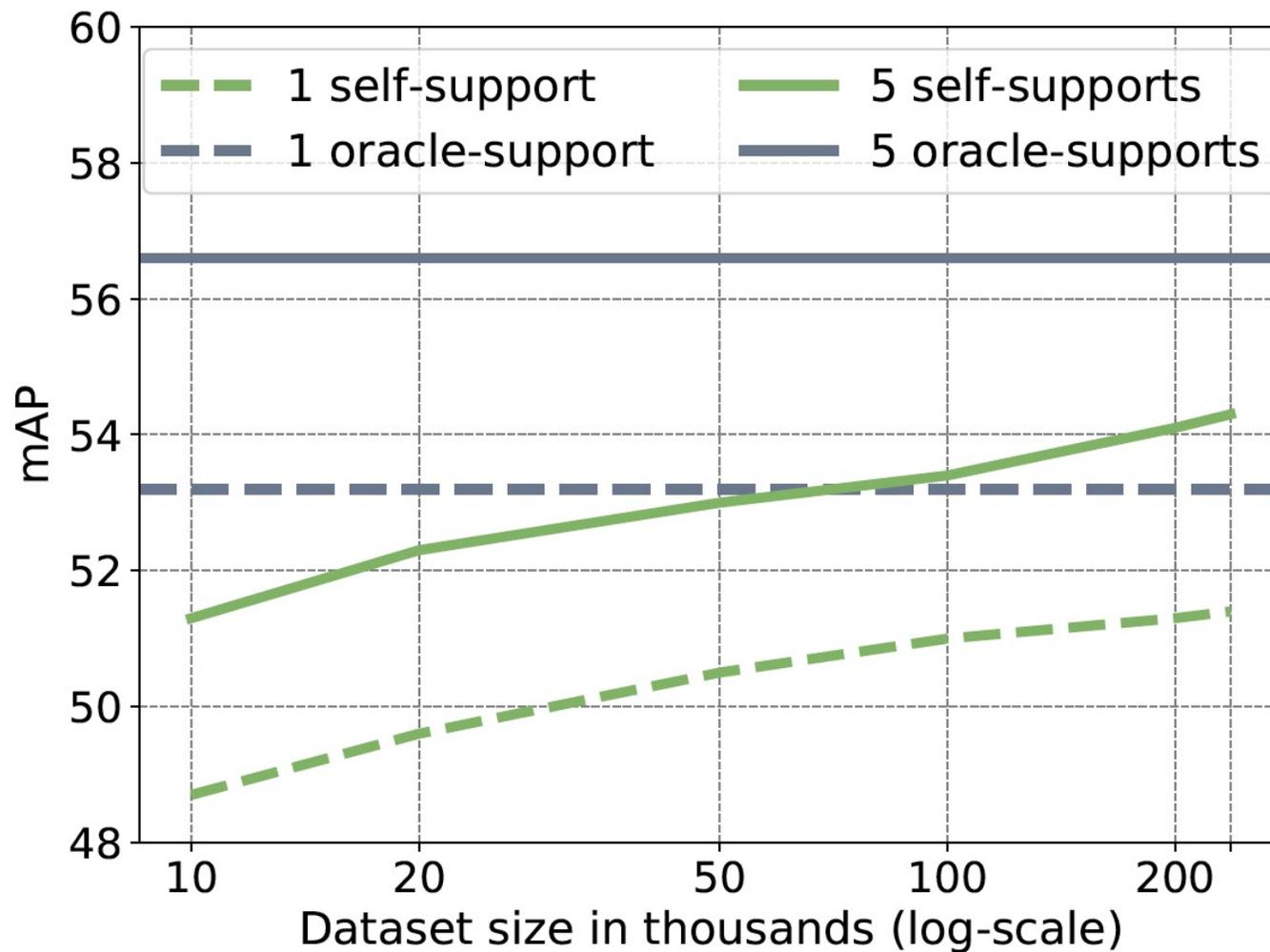
Video instance segmentation results

Self-support examples



Query video

Scalable video instance segmentation?



//. Learning without task assumption



Yingjun Du

University of Amsterdam



Xiantong Zhen

University of Amsterdam



Ling Shao

Inception Institute of AI



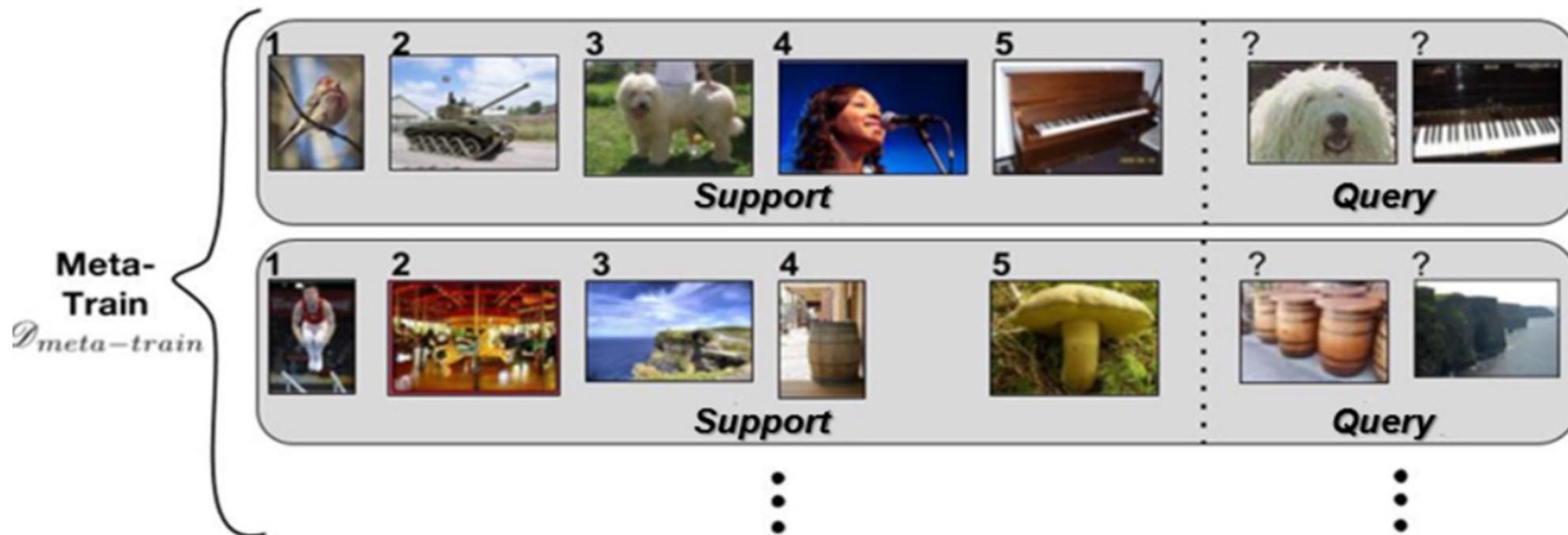
Cees Snoek

University of Amsterdam

MetaNorm: Learning to Normalize Few-Shot Batches Across Domains. In *ICLR* 2021.

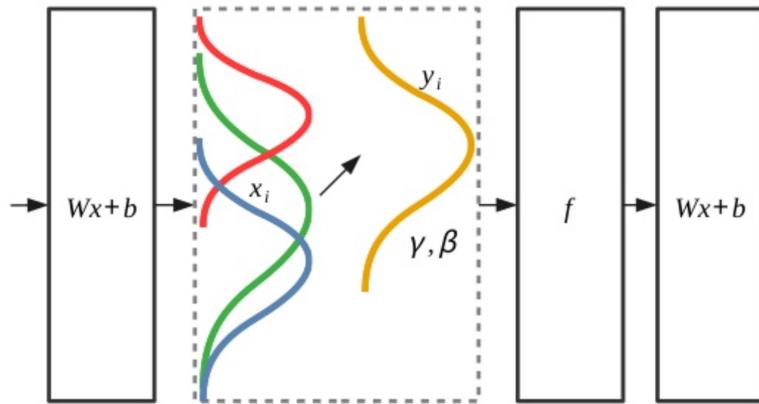
Few-shot meta-learning

5-way, 1-shot



Deep learning work horse: batch normalization

Stabilize the distribution of internal activations during training



$$\mu_{\mathcal{B}} = \frac{1}{M} \sum_{i=1}^M a_i, \quad \sigma_{\mathcal{B}}^2 = \frac{1}{M} \sum_{i=1}^M (a_i - \mu_{\mathcal{B}})^2$$

$$a'_i \leftarrow \mathbf{BN}(a_i) \equiv \gamma \hat{a}_i + \beta, \quad \text{where,} \quad \hat{a}_i = \frac{a_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

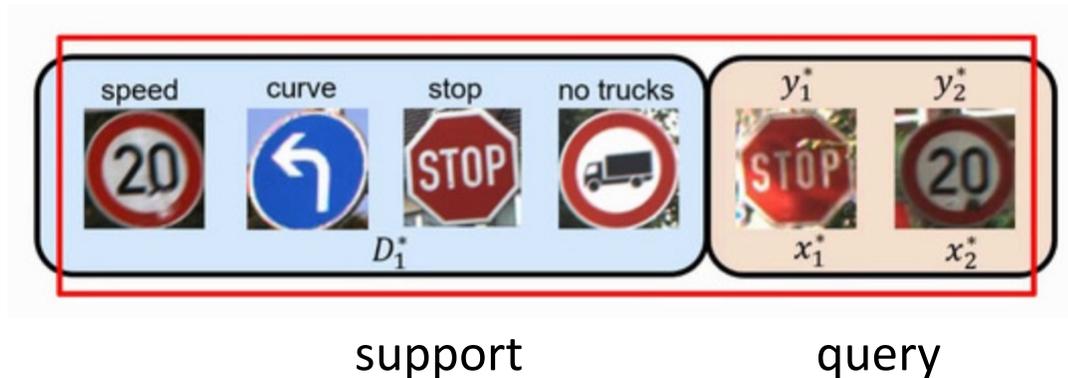
Challenge I: batch statistics become **unstable** with **small batch sizes**

Challenge II: **distribution shift** between source and target domains

Transductive Batch Normalization

Compute batch statistics by using all available query data

Transductive



Requirement to have test set samples available limits real-world use

TaskNorm

Identified the **limiting assumption** of the transductive setting

Leverages statistics from both **layer and instance** normalization

Better than batch norm, sometimes better than transductive.

Our proposal: MetaNorm

Leverage the meta-learning setting

Infer statistics from the support set that better match the query set

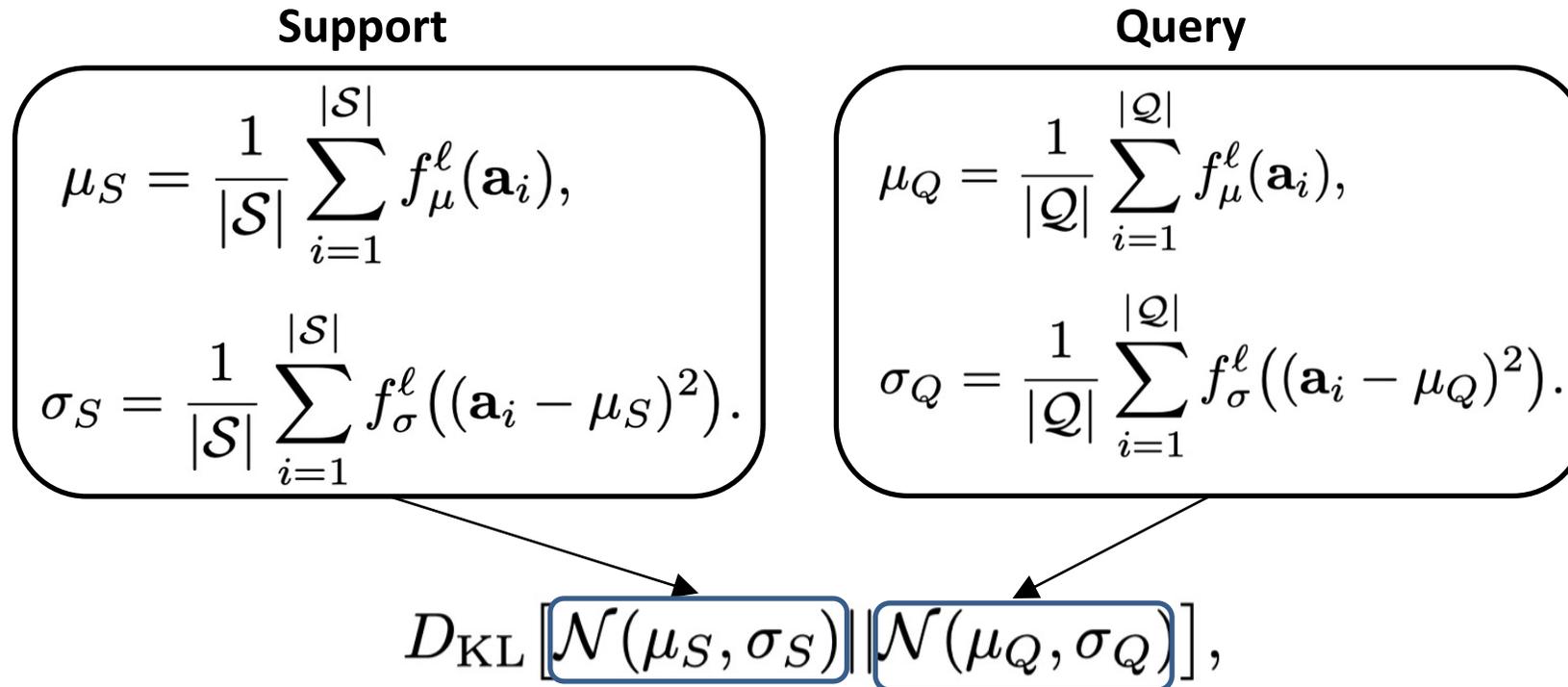
$$D_{\text{KL}} [q(m|S) || p(m|Q)]$$

Distribution inferred from **support set**

Distribution inferred from **query set**

Achieve **adaptive** batch normalization

Meta-training optimization



Hypernetworks $f_{\mu}^{\ell}, f_{\sigma}^{\ell}$, generate (μ_S, σ_S) and (μ_Q, σ_Q) from the support and query sets, for calculating the **KL term** during meta-training optimization.

Meta-testing

Given a test task, the learned hypernetworks $f_\mu^\ell, f_\sigma^\ell$ take the support set as input to **generate** normalization statistics **directly** for the query set.

$$\mu_S = \frac{1}{|S|} \sum_{i=1}^{|S|} f_\mu^\ell(\mathbf{a}_i), \quad \sigma_S = \frac{1}{|S|} \sum_{i=1}^{|S|} f_\sigma^\ell((\mathbf{a}_i - \mu_S)^2).$$

$$a' = \gamma \left(\frac{a - \mu_S}{\sqrt{\sigma_S^2 + \epsilon}} \right) + \beta,$$

Effect of the KL term

	<i>Label gap</i>		<i>Distribution gap</i>				
	Few-shot classification		Domain generalization				
MetaNorm	5-way, 1-shot	5-way, 5-shot	<i>Photo</i>	<i>Art</i>	<i>Cartoon</i>	<i>Sketch</i>	<i>Mean</i>
w/o KL	34.3 \pm 1.5	50.7 \pm 0.8	88.96	71.25	65.37	69.28	73.72
w/ KL	46.8 \pm 1.6	60.1 \pm 0.8	95.99	85.01	78.63	83.17	85.70

Effective for both *few-shot* classification and *many-shot* domain generalization

Comparison with other batch norms

	ProtoNets		MAML	
	5-way, 1-shot	5-way, 5-shot	5-way, 1-shot	5-way, 5-shot
TBN	45.9 \pm 0.6	65.5 \pm 0.9	45.5 \pm 1.8	59.7 \pm 0.9
CBN	47.8 \pm 0.6	66.7 \pm 0.5	20.1 \pm 0.0	20.2 \pm 0.2
TaskNorm	47.5 \pm 0.6	65.3 \pm 0.5	42.0 \pm 1.7	58.1 \pm 0.9
MetaNorm	48.1 \pm 1.6	65.9 \pm 0.9	46.8 \pm 1.6	60.1 \pm 0.8

MetaNorm outperforms transductive and non-transductive normalizations

Few-shot domain generalization

	MAML	
	5-way, 1-shot	5-way, 5-shot
TBN	28.7 \pm 1.8	49.3 \pm 0.8
CBN	20.0 \pm 0.0	20.1 \pm 0.2
TaskNorm	26.9 \pm 1.7	47.4 \pm 0.8
MetaNorm	32.7 \pm 1.7	51.9 \pm 0.9

MetaNorm allows for batch normalization of small batches across domains.

III. Learning without domain assumption



Zehao Xiao

University of Amsterdam



Xiantong Zhen

University of Amsterdam



Ling Shao

Inception Institute of AI

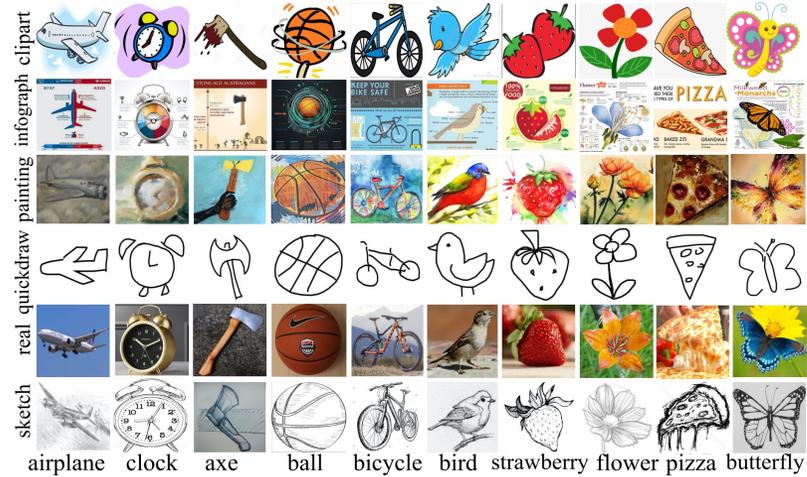


Cees Snoek

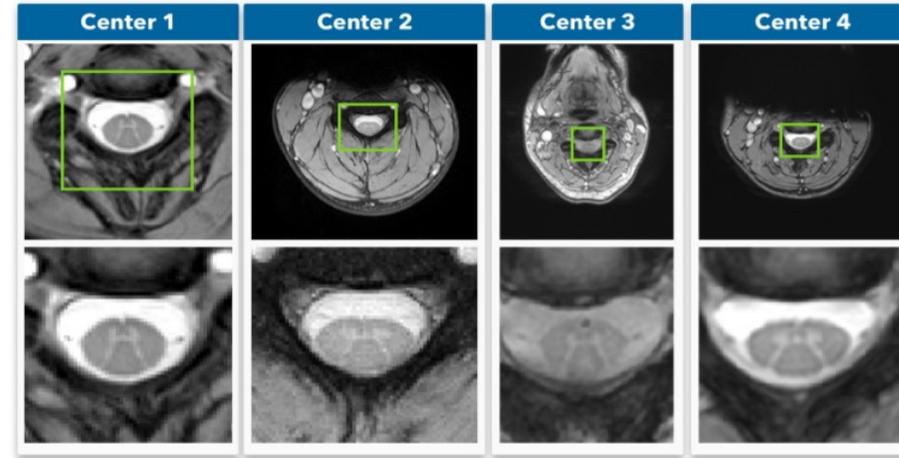
University of Amsterdam

Learning to Generalize across Domains on Single Test Samples. *Submitted.*

Distribution gaps are a fact of life



Images with different style



Medical images from different devices



Daylight



Night



Downtown

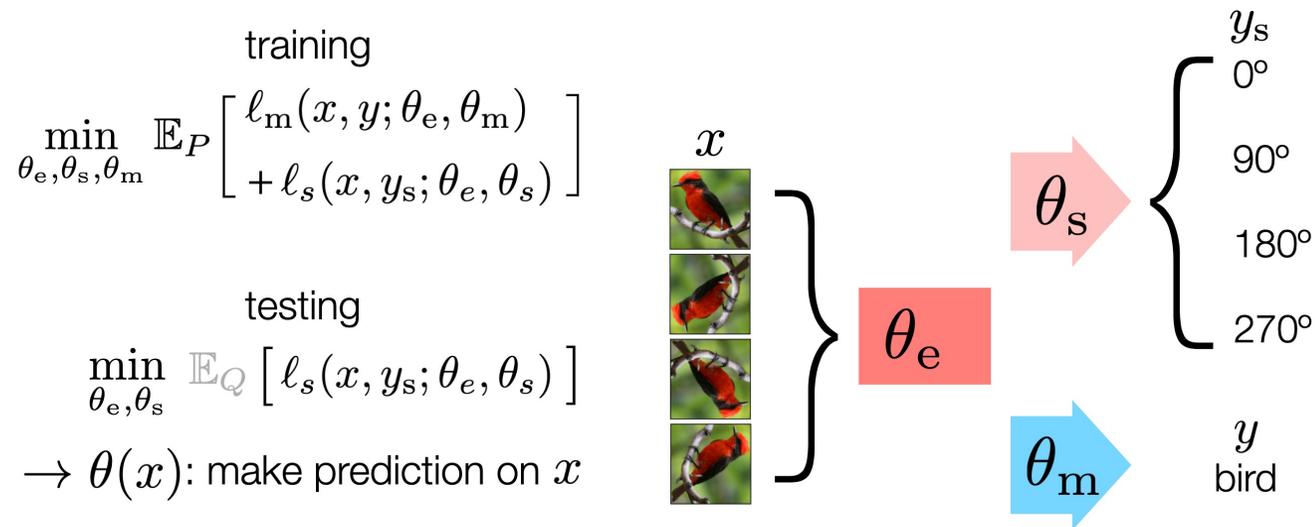


Suburban

Autopilot data in different environments

Test-time training

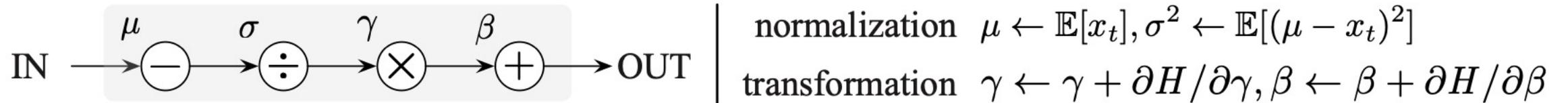
Update model parameters by **self-supervision** before prediction



Needs additional self-supervised **model**, plus **fine-tuning**

Test-time adaptation

Normalize test-batch predictions by entropy minimization

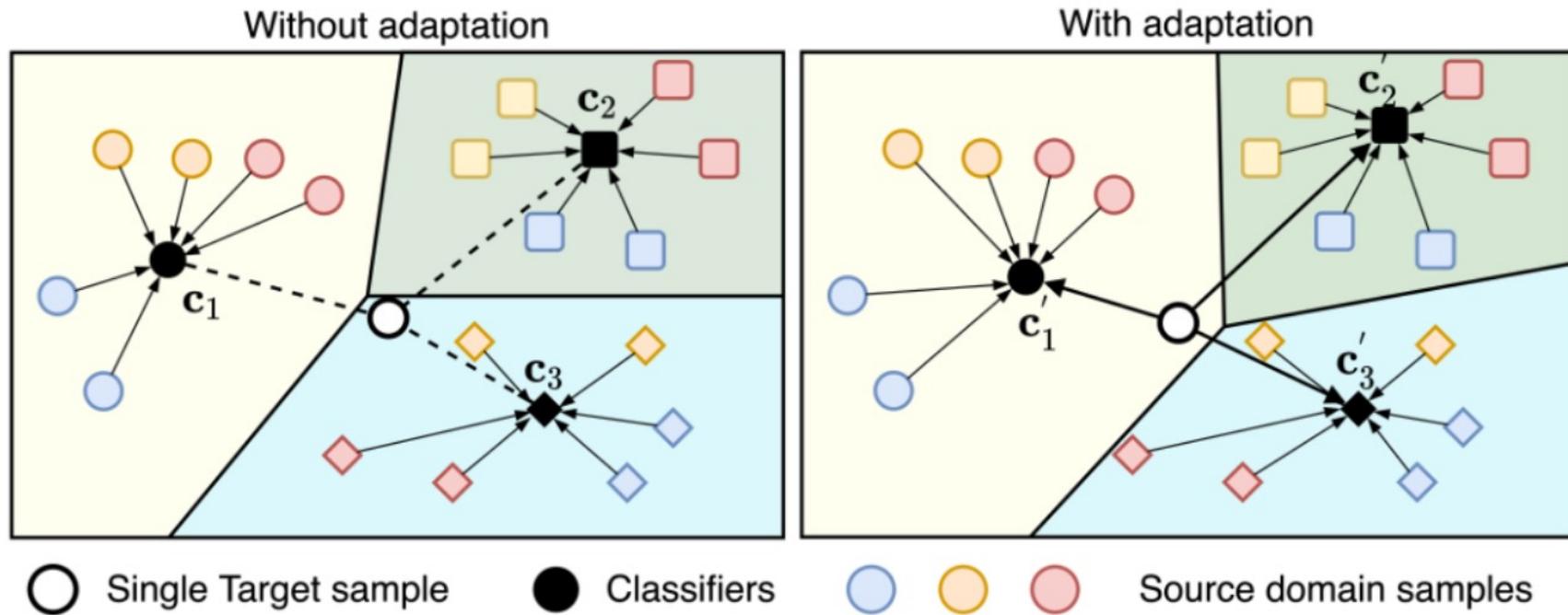


Needs a **batch** to be from the same domain, plus **fine-tuning**

Outperforms test-time training

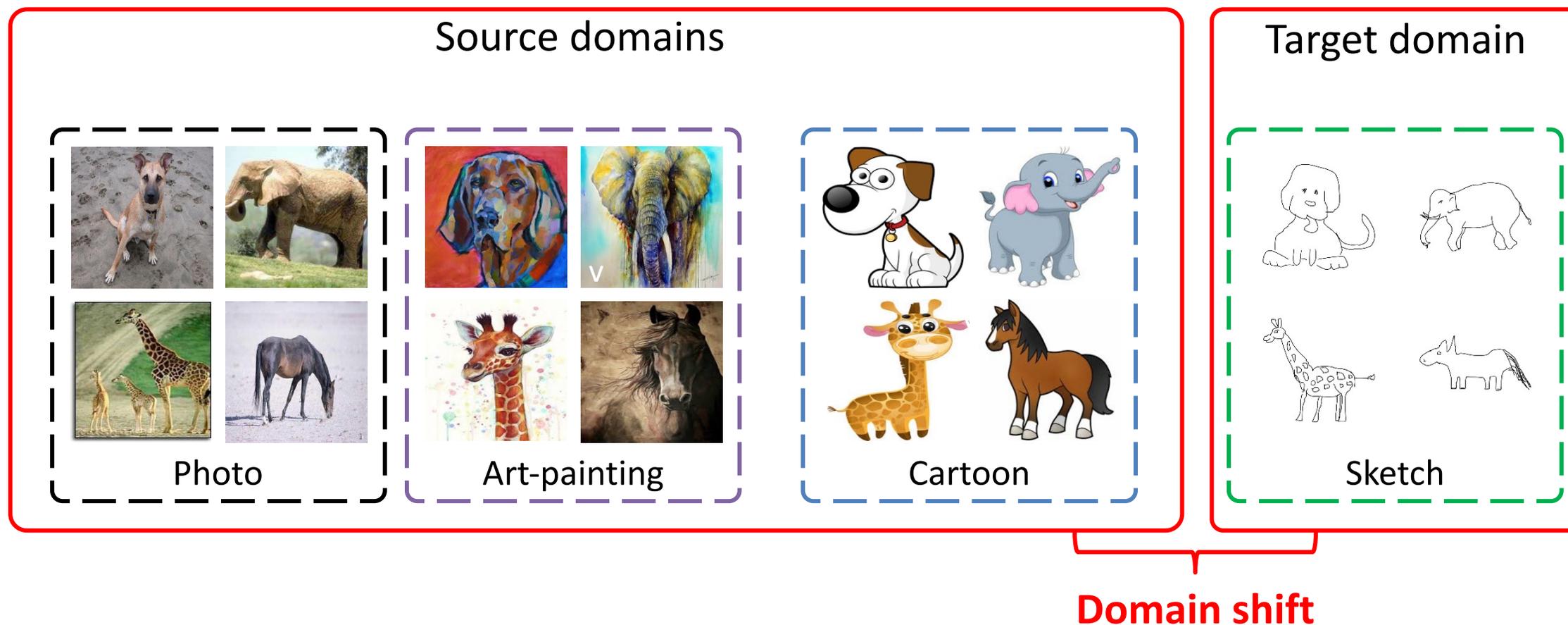
Key idea

Adapt source domain classifiers to each individual target sample



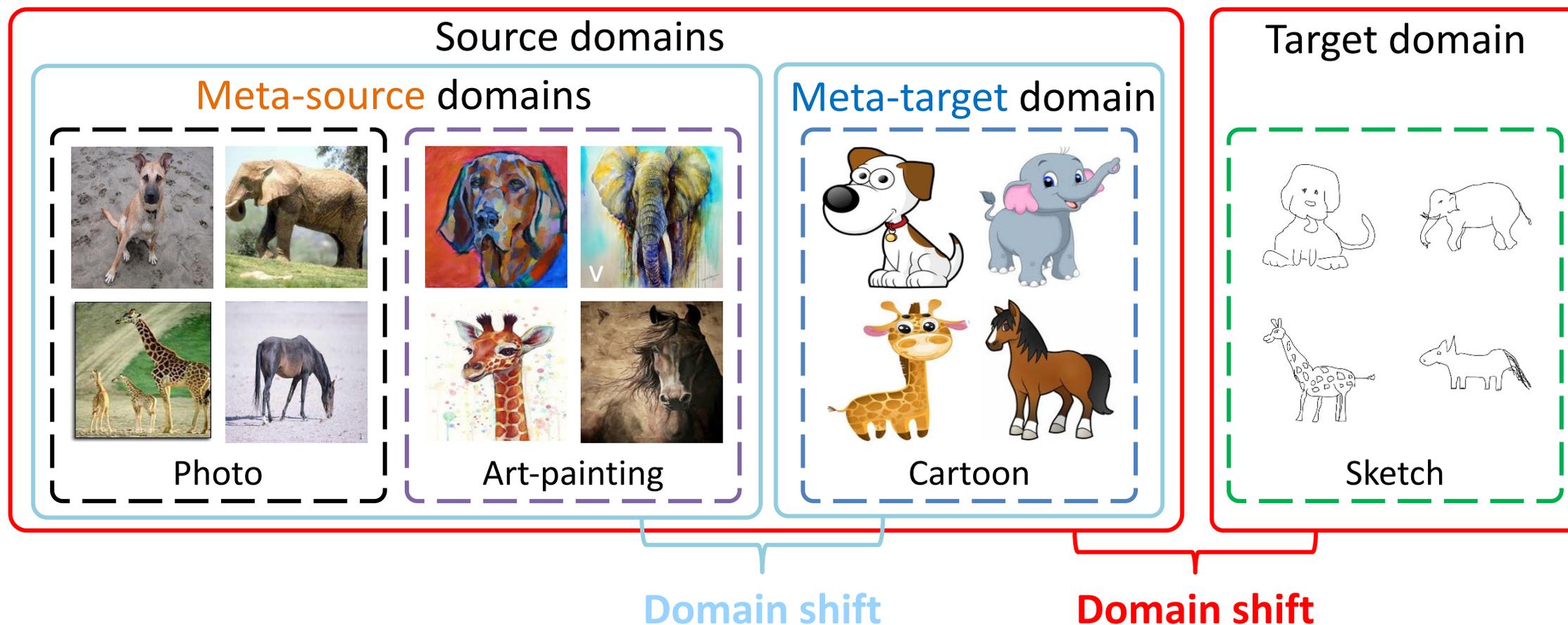
Meta-learning framework

Mimic shift between source and target by shift among source domains



Meta-learning framework

Mimic shift between source and target by shift among source domains



Adaptation as variational inference

Incorporate the **test sample** as a conditional for generating model parameters

$$\log p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \mathcal{T}') = \log \int p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \boldsymbol{\theta}_{t'}) p(\boldsymbol{\theta}_{t'} | \mathcal{T}') d\boldsymbol{\theta}_{t'},$$

Meta-target

Adaptation as variational inference

Incorporate the **test sample** as a conditional for generating model parameters

$$\log p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \mathcal{T}') = \log \int p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \boldsymbol{\theta}_{t'}) p(\boldsymbol{\theta}_{t'} | \mathcal{T}') d\boldsymbol{\theta}_{t'},$$

Meta-target

Intractable during inference, so we approximate by source domain similarity

$$\geq \mathbb{E}_{q(\boldsymbol{\theta}_{t'})} [\log p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \boldsymbol{\theta}_{t'})] - \mathbb{D}_{\text{KL}}[q(\boldsymbol{\theta}_{t'} | \mathbf{x}_{t'}, \mathcal{S}') || p(\boldsymbol{\theta}_{t'} | \mathcal{T}')].$$

Meta-source

Adaptation as variational inference

Incorporate the **test sample** as a conditional for generating model parameters

$$\log p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \mathcal{T}') = \log \int p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \boldsymbol{\theta}_{t'}) p(\boldsymbol{\theta}_{t'} | \mathcal{T}') d\boldsymbol{\theta}_{t'},$$

Meta-target

Intractable during inference, so we approximate by source domain similarity

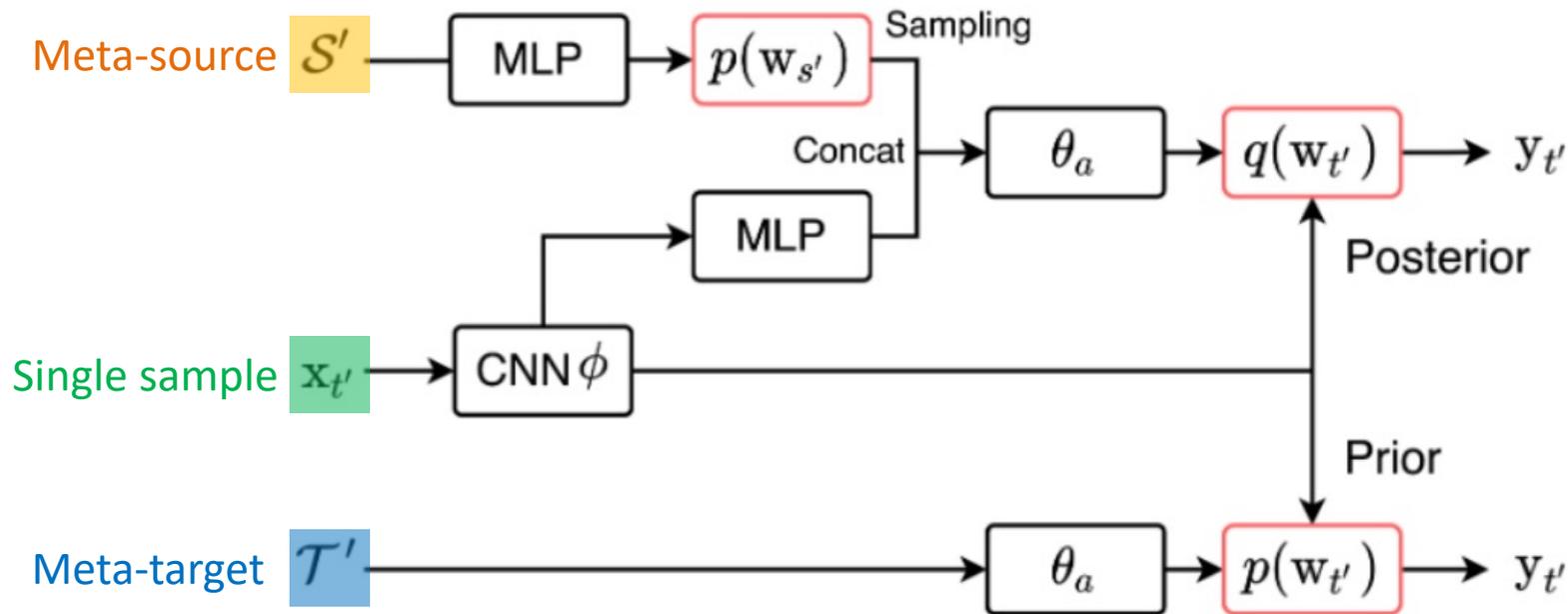
$$\geq \mathbb{E}_{q(\boldsymbol{\theta}_{t'})} [\log p(\mathbf{y}_{t'} | \mathbf{x}_{t'}, \boldsymbol{\theta}_{t'})] - \mathbb{D}_{\text{KL}}[q(\boldsymbol{\theta}_{t'} | \mathbf{x}_{t'}, \mathcal{S}') || p(\boldsymbol{\theta}_{t'} | \mathcal{T}')].$$

Meta-source

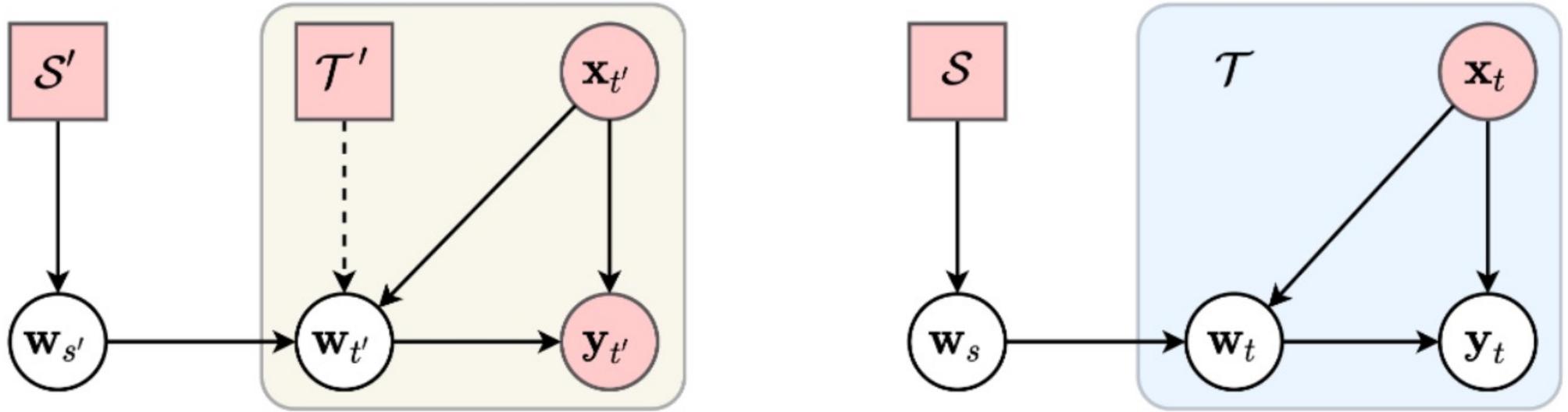
Our model **learns the ability to adapt** the **meta-source** model to each **meta-target** instance across different domain shifts

Computational feasibility

We divide the model θ into a feature extractor ϕ and a classifier w .
 ϕ is shared across domains, while w is trained to be adapted



Generalization at test-time

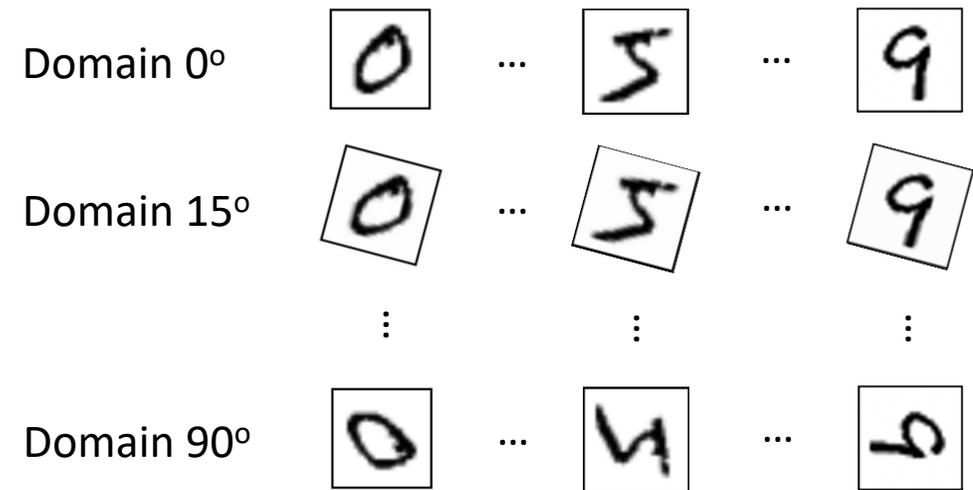


(a) Training on meta-source (S') and meta-target (T') domains (b) Testing on the unseen domain (T)

Adaptation is achieved by generating w_t for each target sample with only **one forward pass** using an amortization inference network

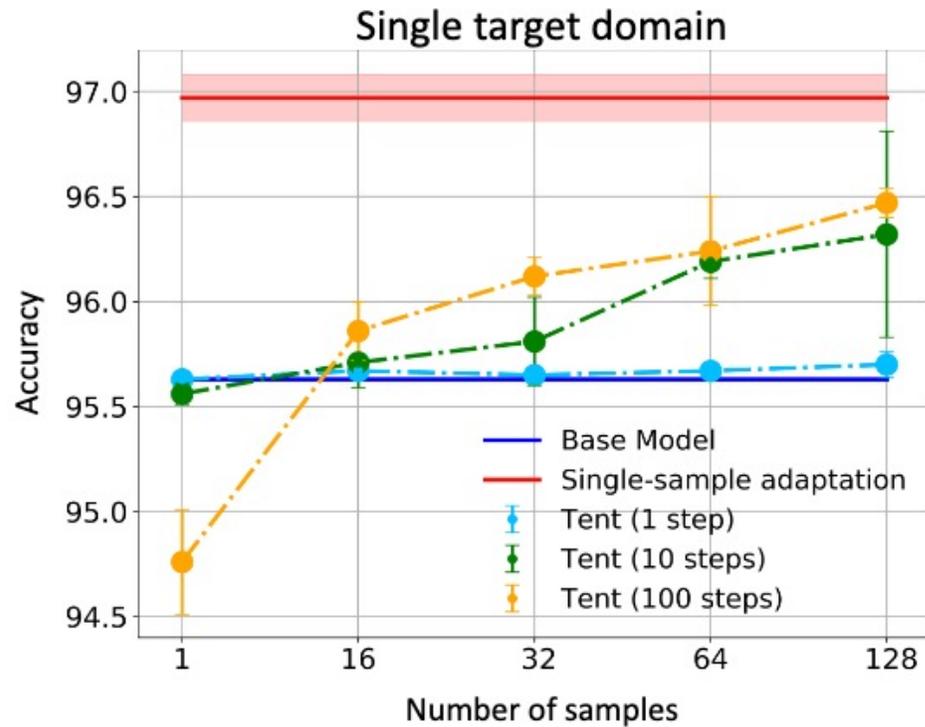
Comparison with test-time adaptation (Tent)

Results on Rotated-MNIST

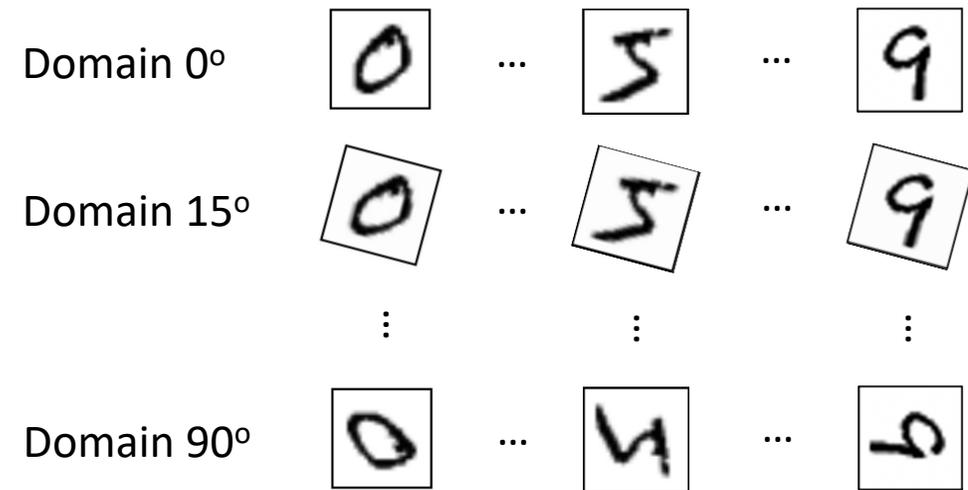


Comparison with test-time adaptation (Tent)

Tent works well with a **large batch** of samples from a **single target** domain



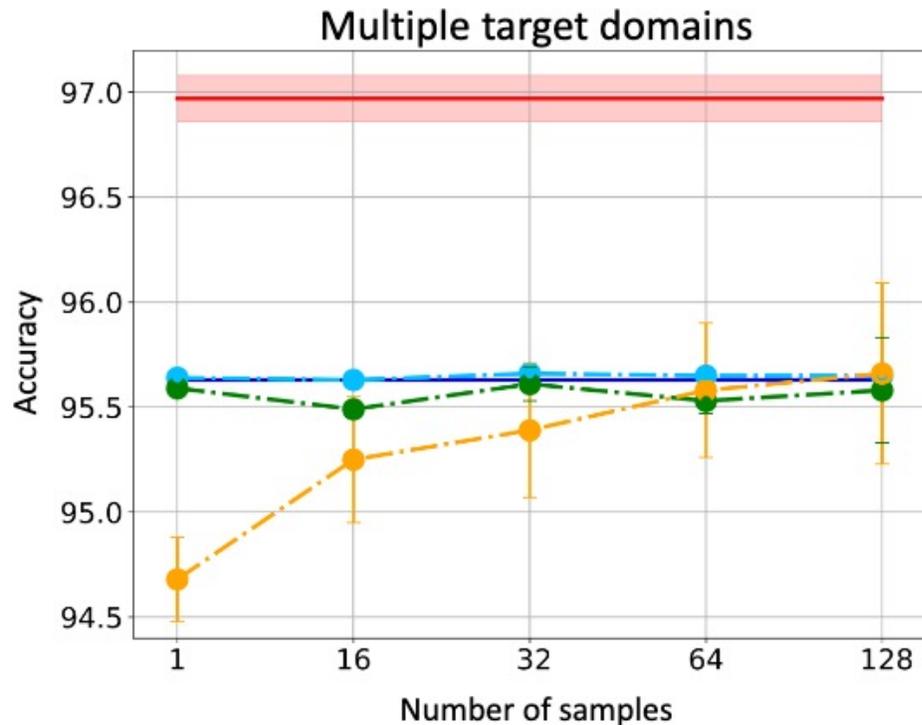
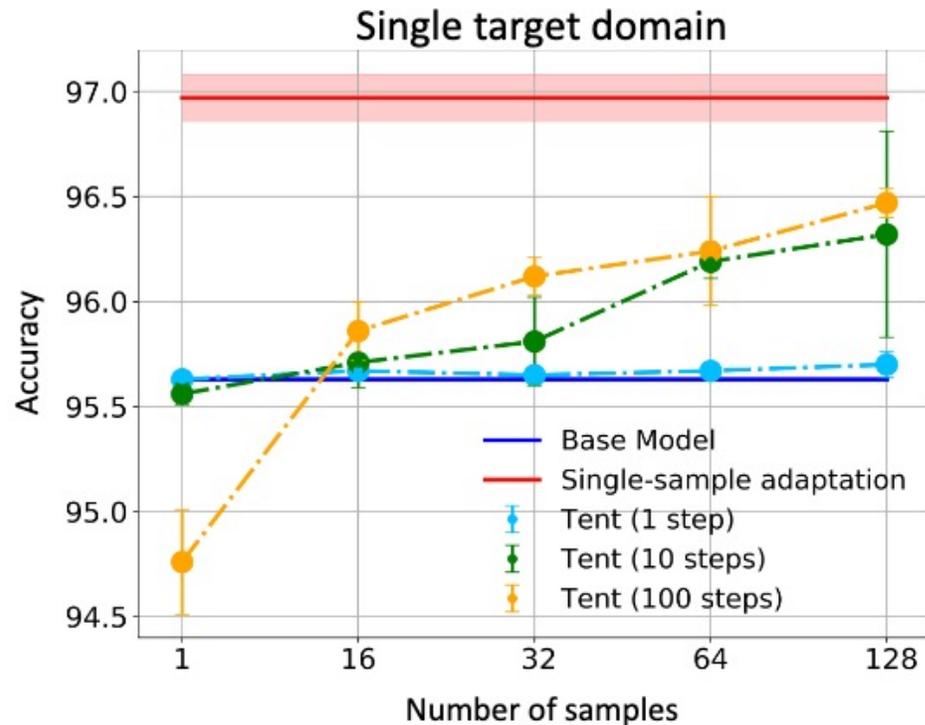
Results on Rotated-MNIST



Comparison with test-time adaptation (Tent)

Tent works well with a **large batch** of samples from a **single target** domain

We outperform with a **single sample**, especially for **multiple target** domains



More comparisons

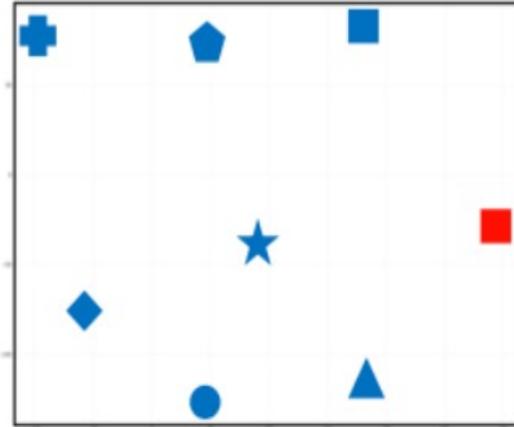
The better the base network, the more we gain.

	PACS benchmark		Office-Home benchmark		
	ResNet-18	ResNet-50	ResNet-18	ResNet-50	
Wang et al. ICLR 2021	83.09	86.23	64.13	67.99	} Test-time adaptation
Dubey et al. CVPR 2021	--	84.50		68.90	
Zhou et al. ECCV 2020	83.70	84.90	65.63	67.66	Domain generation
Seo et al. ECCV 2020	85.11	86.64	62.90	--	Normalization
<i>Ours</i>	84.15	87.51	66.02	71.07	

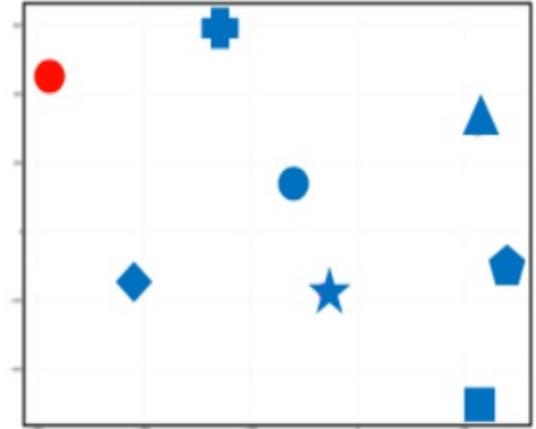
Failure cases



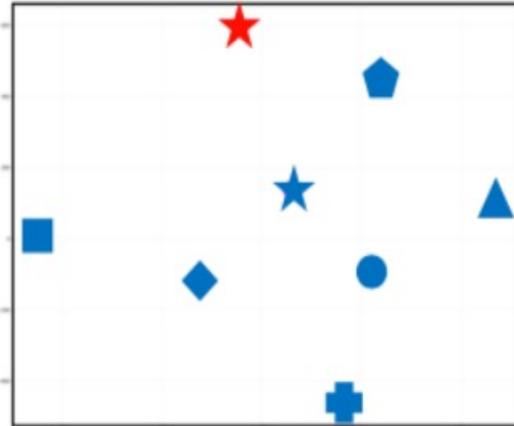
Label: Guitar
Prediction: Person



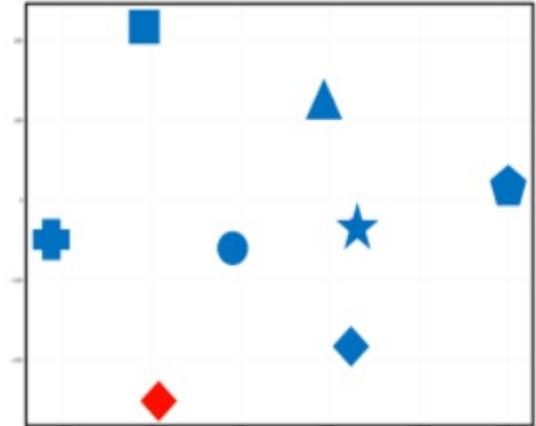
Label: Horse
Prediction: House



Label: Dog
Prediction: Giraffe



Label: Elephant
Prediction: Horse



Classifiers:

◆ Elephant

★ Dog

⬠ Giraffe

■ Guitar

● Horse

⊕ House

▲ Person

Conclusions

Need for **real-world learning** across domains, labels, tasks and with fairness.

Need to **question** common learning **assumptions**.

Label, task and **domain** assumptions can be relaxed during learning.

Thank you

Contact info

Prof. dr. Cees Snoek

<https://www.ceessnoek.info>

cgmsnoek@uva.nl

@cgmsnoek