

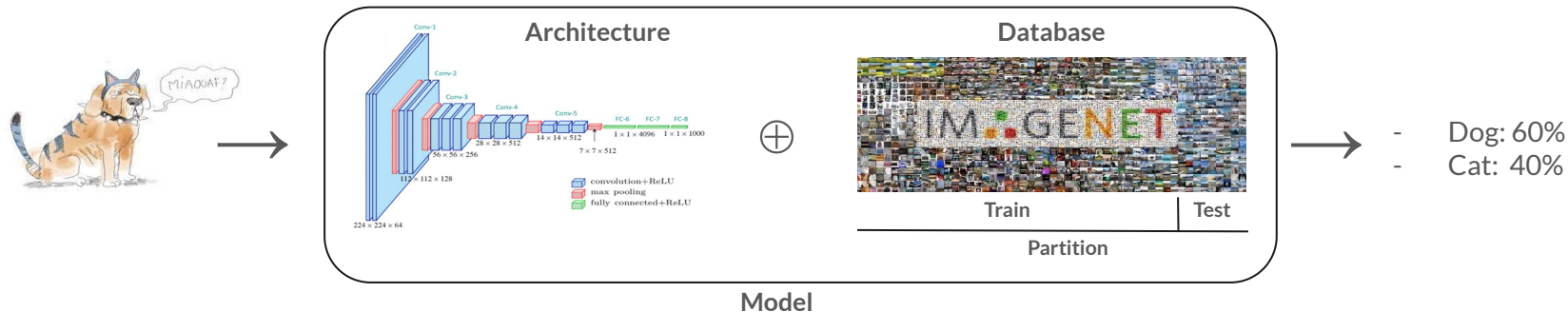


# Consistent K-fold cross validation:

Application to image/video object detection datasets

C. Barrelet, M. Chaumont, G. Subsol, V. Creuze, and M. Gouttefarde

# What's a K-fold?



| Model | Folds: | Train | Test | $\theta$      |
|-------|--------|-------|------|---------------|
| $M_1$ |        |       |      | $\theta(M_1)$ |
| $M_2$ |        |       |      | $\theta(M_2)$ |
| $M_3$ |        |       |      | $\theta(M_3)$ |
| $M_4$ |        |       |      | $\theta(M_4)$ |
| $M_5$ |        |       |      | $\theta(M_5)$ |

Architecture performance:

$$\theta(\mathcal{M}) \approx \bar{\theta}_M = \sum_{i=1}^K \theta(M_i) / K$$

Performance sensitivity:

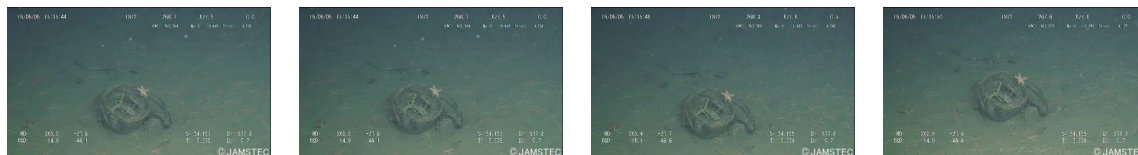
$$\sigma^2(\mathcal{M}) \approx \sigma_M^2 = \sum_{i=1}^K (\theta(M_i) - \bar{\theta}_M)^2 / K$$

# Why a K-fold?

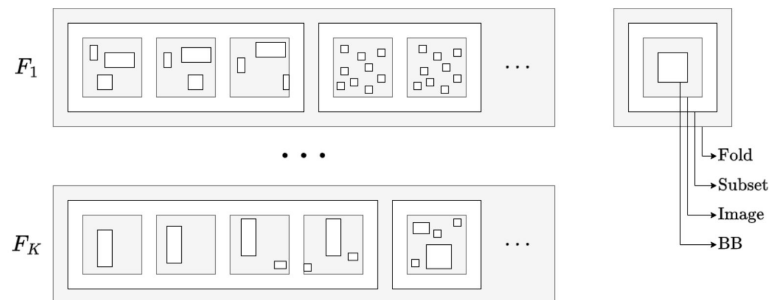


- Not necessary for “large” reference datasets (ImageNet, COCO, etc.)
  - Law of large numbers
- Necessary for “smaller” datasets
  - Law of large numbers cannot hold anymore
- Easier architectures comparison

**But, what if the dataset can be subdivided into subsets of images sharing similar features?**



## What does it mean for object detection datasets?



# Internal variability you say?

## K matters!

- Moderate  $K$  values (10-20)  $\rightarrow \sigma_M^2$
- Low  $K$  values (2-5)  $\rightarrow \sigma_M^2$

## Size too!

- Dataset size  $N$   $\rightarrow \sigma_M^2$
- Better use  $K=5$  or  $K=10$  than  $K=N$

## There is a performance-variance tradeoff...

- Estimator must maximize the performance and minimize the variance

## Random induces internal variability.

## Maybe smooth the results with J-K-folds?

- How to choose J and K?
- Dataset more heterogeneous  $\rightarrow$  More repetitions...
- More complexity...

## What about saving complexity by ensuring internal consistency?

Ron Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2, San Francisco, CA, USA, 1995, IJCAI'95, p. 1137-1143, Morgan Kaufmann Publishers Inc

Juan D. Rodriguez, Aritz Perez, and Jose A. Lozano, "Sensitivity analysis of k-fold cross validation in prediction error estimation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 32, no. 3, pp.569-575, 2010

Ron Kohavi and David H. Wolpert, "Bias plus variance decomposition for zero-one loss functions," in International Conference on Machine Learning, 1996

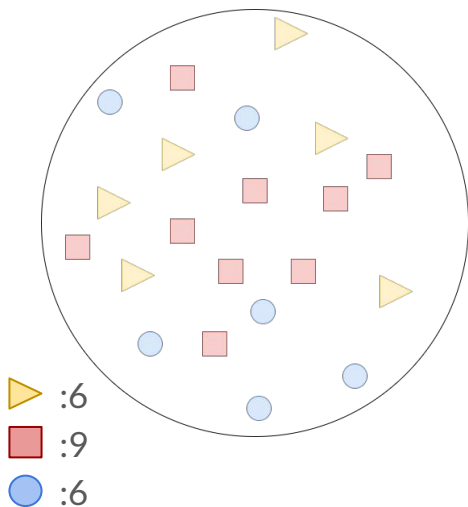
Gaoxia Jiang and Wenjian Wang, "Error estimation based on variance analysis of k-fold cross-validation," Pattern Recognition, vol. 69, pp. 94-106, 2017

Tzu-Tsung Wong and Po-Yang Yeh, "Reliable accuracy estimates from k-fold cross validation," IEEE Transactions on Knowledge and Data Engineering, vol. 32, no. 8, pp. 1586-1594, 2020.

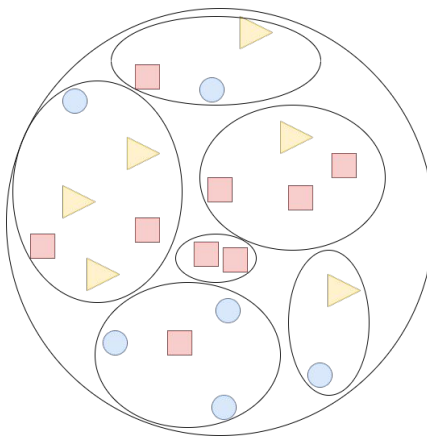
Henry Moss, David Leslie, and Paul Rayson, "Using J-K-fold cross validation to reduce variance when tuning NLP models," in Proceedings of the 27th International Conference on Computational Linguistics, Santa Fe, New Mexico, USA, Aug. 2018, pp. 2978-2989, Association for Computational Linguistics.

# Let's ensure consistency!

Stratified-KCV

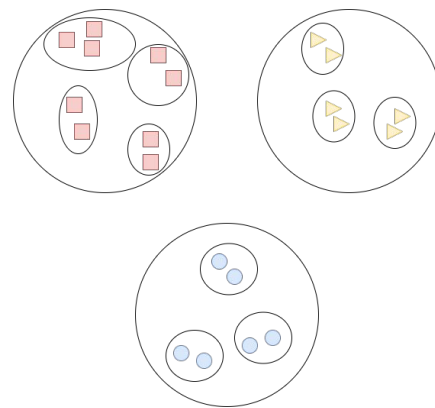


Distributed-KCV



Predetermined and hand-crafted features...

Distributed-Stratified-KCV



## Classification task only

# A general method to build consistent KCV

## Definition of descriptors at BB and images levels

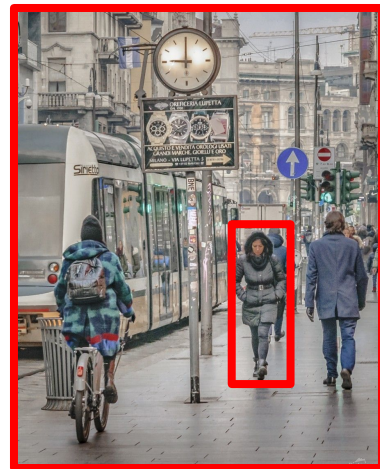
- **BB-level:**
  - a class
  - features describing geometrical and intensity properties of the **thumbnail**
- **Image-level:**
  - features describing geometrical and intensity properties of the **image**

## An histogram-based framework

- Consider an image dataset of  $N$  subsets identified by the index  $s \in \{1, \dots, N\}$
- Assume each subset  $s$  characterized by a set of  $m$  histograms denoted  $\mathbf{H}_h^s$  with  $h \in \{0, \dots, m\}$

| k = 1 |       |        |        |
|-------|-------|--------|--------|
| s = 1 | s = 7 | s = 18 | s = 21 |
| $H_0$ | $H_0$ | $H_0$  | $H_0$  |
| ...   | ...   | ...    | ...    |
| $H_m$ | $H_m$ | $H_m$  | $H_m$  |

| k = 2 |        |        |        |
|-------|--------|--------|--------|
| s = 2 | s = 11 | s = 17 | s = 34 |
| $H_0$ | $H_0$  | $H_0$  | $H_0$  |
| ...   | ...    | ...    | ...    |
| $H_m$ | $H_m$  | $H_m$  | $H_m$  |



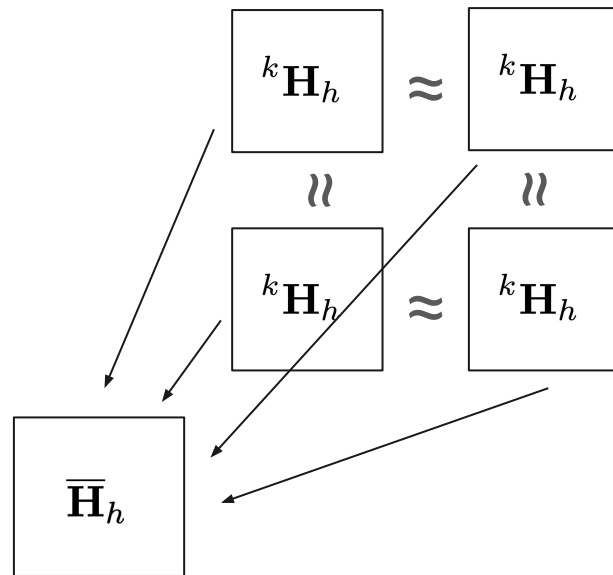
# A general method to build consistent KCV

| k = 1          |                |                |                |
|----------------|----------------|----------------|----------------|
| s = 1          | s = 7          | s = 18         | s = 21         |
| H <sub>0</sub> | H <sub>0</sub> | H <sub>0</sub> | H <sub>0</sub> |
| ...            | ...            | ...            | ...            |
| H <sub>m</sub> | H <sub>m</sub> | H <sub>m</sub> | H <sub>m</sub> |

The average of the  $h^{\text{th}}$  histogram for the fold  $k \rightarrow {}^k\mathbf{H}_h = \frac{1}{|{}^k\mathcal{S}|} \sum_{s \in {}^k\mathcal{S}} \mathbf{H}_h^s$

The desired list of histograms  $\rightarrow \bar{\mathbf{H}}_h = \frac{1}{N} \sum_{s=1}^N \mathbf{H}_h^s$

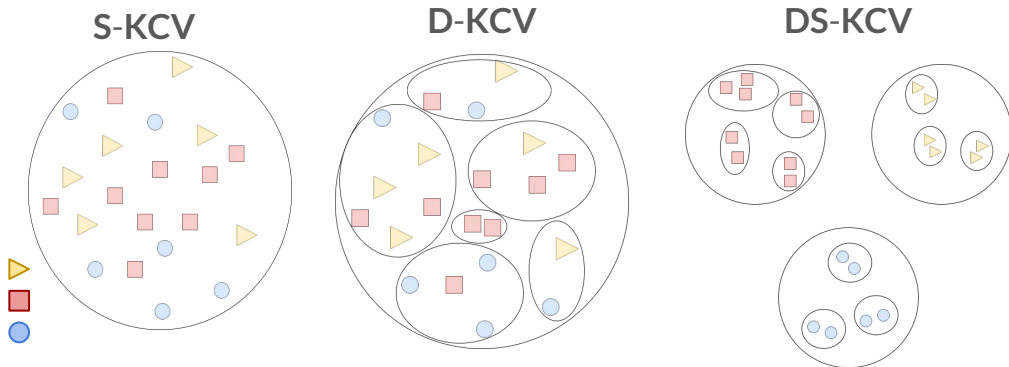
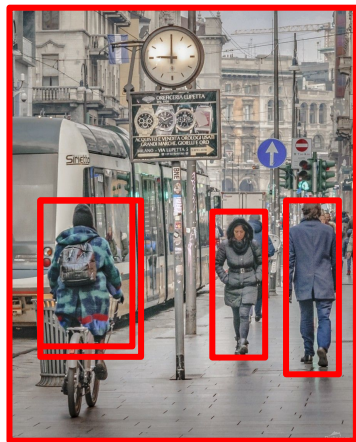
$$\mathbf{a}^* = \underset{\mathbf{a} \in \{1, \dots, K\}^N}{\operatorname{argmin}} \left( \sum_{k=1}^K \sum_{h=0}^m D \left( {}^k\mathbf{H}_h, \bar{\mathbf{H}}_h \right) \right)$$





# A general method to build consistent KCV

3 different descriptors:



$$\mathbf{a}^* = \underset{\mathbf{a} \in \{1, \dots, K\}^N}{\operatorname{argmin}} \left( \sum_{k=1}^K \alpha_0 D \left( \overset{\text{S-KCV}}{\overset{\uparrow}{\mathbf{H}_0^k}}, \bar{\mathbf{H}}_0 \right) + \alpha_1 D \left( \overset{\text{D-KCV}}{\overset{\uparrow}{\mathbf{H}_1^k}}, \bar{\mathbf{H}}_1 \right) + \alpha_2 D \left( \overset{\text{DS-KCV}}{\overset{\uparrow}{\mathbf{H}_2^k}}, \bar{\mathbf{H}}_2 \right) \right)$$

# A greedy optimization algorithm

- There are  $K^N$  potential assignments
- Similar to the *Multi-way number partitioning problem* (NP-hard)

---

**Algorithm 1:** A greedy algorithm to build a balanced  $K$ -fold

---

**Inputs :**  $K, N, \{\forall h, \bar{\mathbf{H}}_h\}$   
**Output:**  $\mathbf{a}$  ▷ Videos assignation

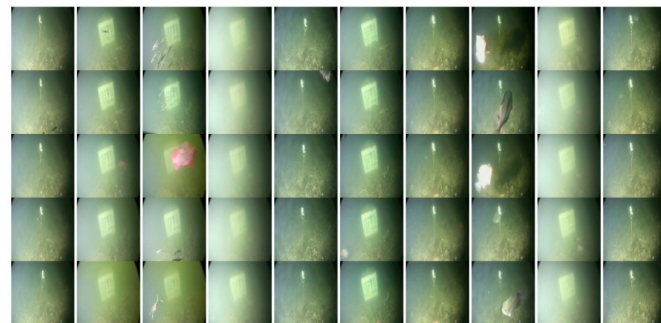
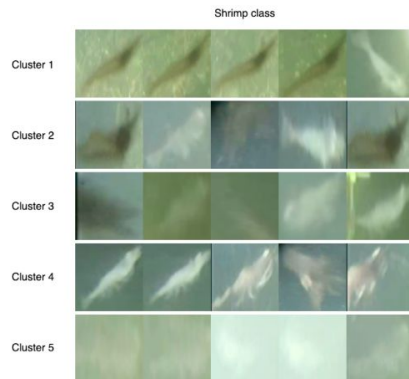
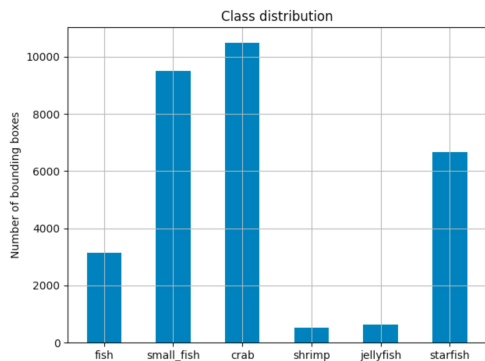
1  $\mathbf{a} \leftarrow [0, 0, \dots, 0]$  ▷ Initialization  
2  $\mathcal{V} \leftarrow \{0, \dots, N\}$   
/\* While there is a video to dispatch among the folds \*/  
3 **while**  $\mathcal{V} \neq \emptyset$  **do**  
4   **foreach**  $k$  randomly picked from  $\{1, \dots, K\}$  **do**  
5     /\* Find the best video to assign to the fold  $k$  \*/  
6      $\mathbf{e} \leftarrow [0, 0, \dots, 0]$  ▷ Error vector  
7     **forall**  $j$  in  $\mathcal{V}$  **do**  
8        ${}^k\mathcal{V}^{(temp)} \leftarrow {}^k\mathcal{V} \cup j$   
9        $\forall h, {}^k\mathbf{H}_h^{(temp)} \leftarrow \frac{1}{|{}^k\mathcal{V}^{(temp)}|} \sum_{i \in {}^k\mathcal{V}^{(temp)}} H_h^{(i)}$   
10        $\mathbf{e}[j] \leftarrow \sum_{h=0}^m \alpha_h \cdot \text{dist}({}^k\bar{\mathbf{H}}_h^{(temp)}, \bar{\mathbf{H}}_h)$   
11        $j^* \leftarrow \underset{j}{\operatorname{argmin}}(\mathbf{e})$  ▷ Best video to assign  
12        $\mathbf{a}[j^*] \leftarrow k$  ▷ Assign it to the fold  
13        $\mathcal{V} \leftarrow \mathcal{V} \setminus j^*$  ▷ Remove it from  $\mathcal{V}$

---

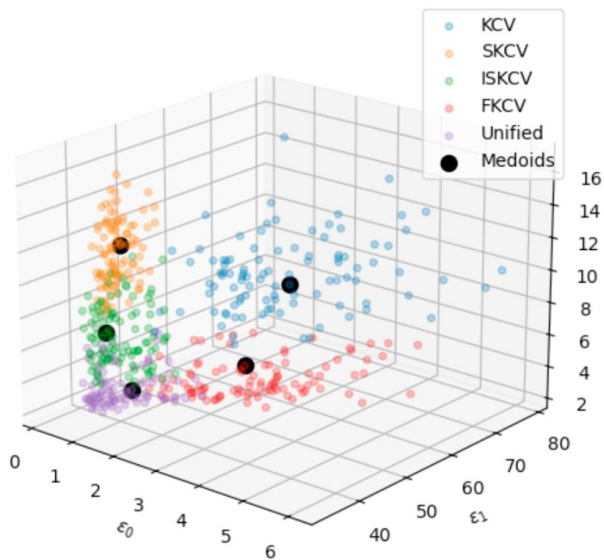
# Application to image and video datasets

Scenario I: A video dataset with class imbalance (**Brackish**)

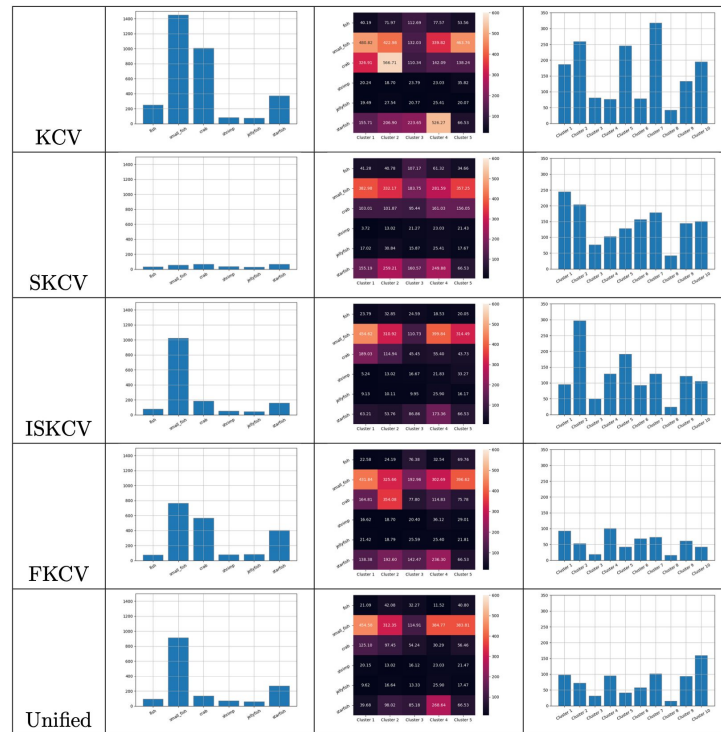
$$\mathbf{a}^* = \underset{\mathbf{a} \in \{1, \dots, K\}^N}{\operatorname{argmin}} \left( \sum_{k=1}^K \alpha_0 D \left( {}^k\mathbf{H}_0, \bar{\mathbf{H}}_0 \right) + \alpha_1 D \left( {}^k\mathbf{H}_1, \bar{\mathbf{H}}_1 \right) + \alpha_2 D \left( {}^k\mathbf{H}_2, \bar{\mathbf{H}}_2 \right) \right)$$



# Consistency results

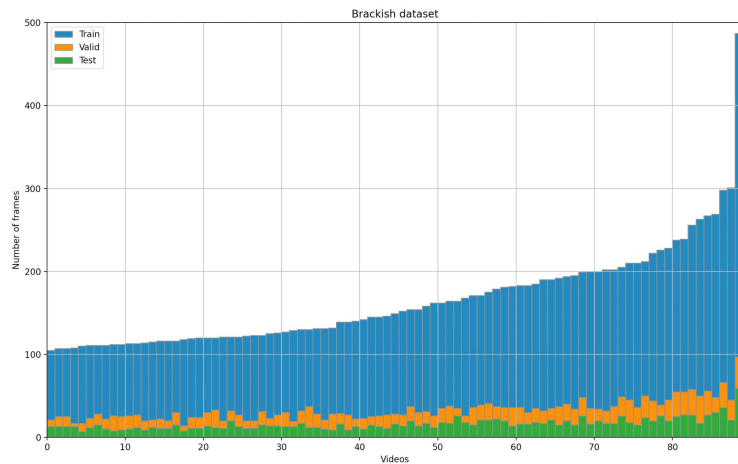


$$\mathbf{a}^* = \underset{\mathbf{a} \in \{1, \dots, K\}^N}{\operatorname{argmin}} \left( \sum_{k=1}^K \alpha_0 D \left( {}^k \mathbf{H}_0, \bar{\mathbf{H}}_0 \right) + \alpha_1 D \left( {}^k \mathbf{H}_1, \bar{\mathbf{H}}_1 \right) + \alpha_2 D \left( {}^k \mathbf{H}_2, \bar{\mathbf{H}}_2 \right) \right)$$



# Results

| Partitioning  | P (%)                              | R (%)                              | F1 (%)                             | mAP@.5 (%)                         | mAP@.5:.95(%)                      |
|---------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| *Biased split | 98.10                              | 98.10                              | 98.10                              | 99.10                              | 83.60                              |
| KCV           | 63.24 $\pm$ 9.61                   | 43.90 $\pm$ 6.75                   | 48.64 $\pm$ 5.74                   | 47.83 $\pm$ 7.16                   | 29.39 $\pm$ 6.35                   |
| SKCV          | 62.76 $\pm$ 9.58                   | 44.55 $\pm$ 6.43                   | 49.47 $\pm$ 6.26                   | 48.97 $\pm$ 6.46                   | 28.54 $\pm$ 3.62                   |
| ISKCV         | 64.06 $\pm$ 10.01                  | 46.02 $\pm$ 6.32                   | 50.92 $\pm$ 7.19                   | <b>50.25 <math>\pm</math> 6.14</b> | <b>29.65 <math>\pm</math> 3.73</b> |
| FKCV          | 62.94 $\pm$ <b>7.17</b>            | 44.95 $\pm$ 5.33                   | 49.15 $\pm$ 5.30                   | 48.89 $\pm$ <b>5.54</b>            | 29.43 $\pm$ 4.00                   |
| Unified       | <b>66.38 <math>\pm</math> 8.07</b> | <b>46.31 <math>\pm</math> 5.22</b> | <b>51.53 <math>\pm</math> 5.29</b> | 50.12 $\pm$ 5.98                   | 29.54 $\pm$ <b>3.52</b>            |



# Application to image and video datasets

Scenario II: An image dataset with class imbalance (PVOC10)

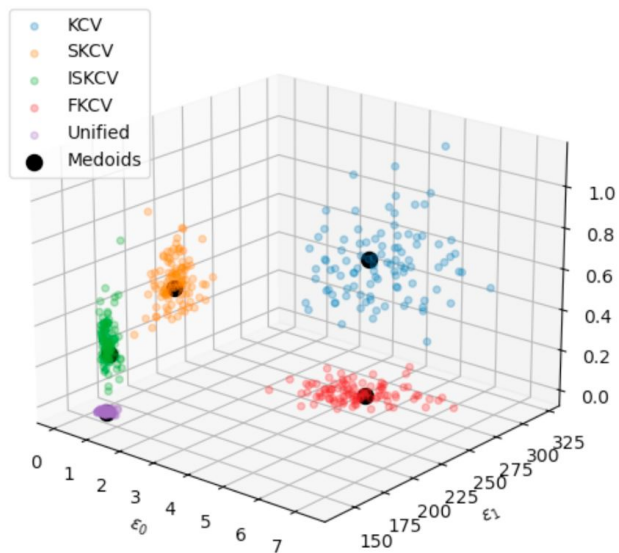


# Conclusion

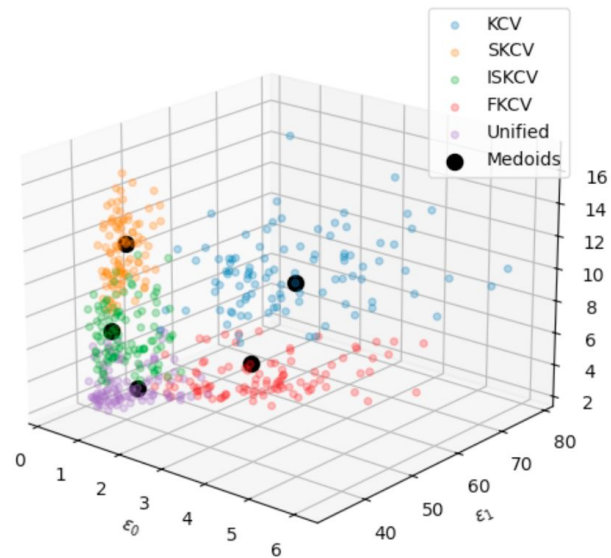


- New paradigm for creating consistent K-fold
- Application to object detection
- Significant results

# Consistency results



PVOC10

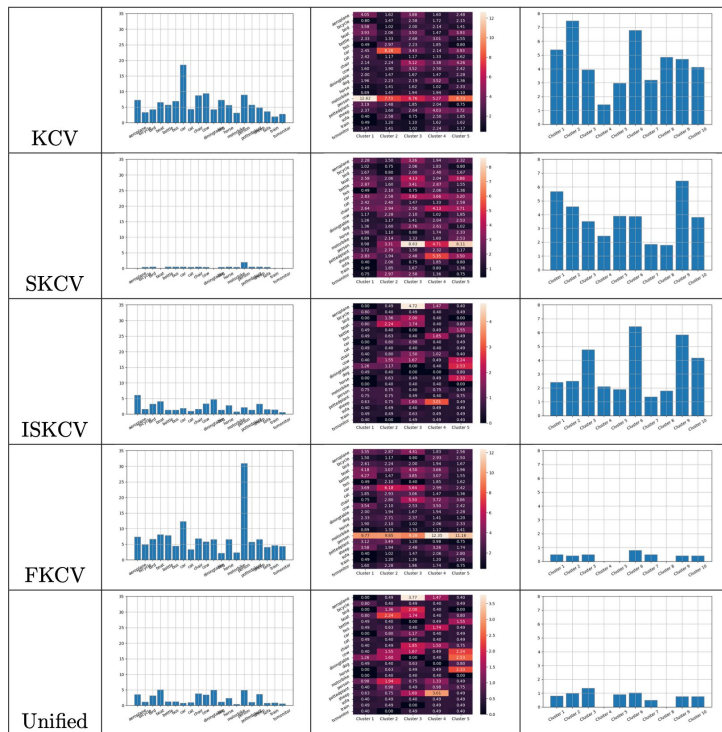


Brackish

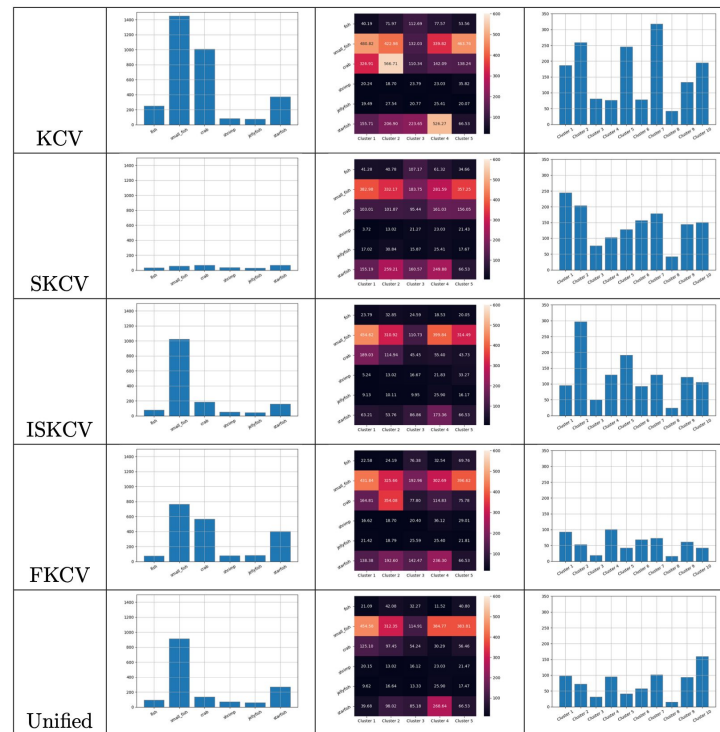


# Consistency results

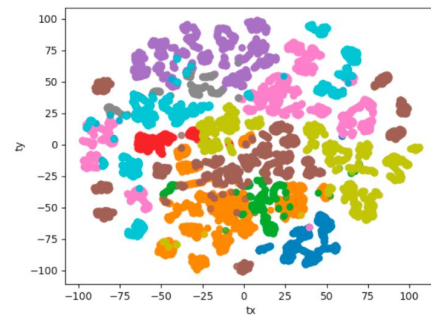
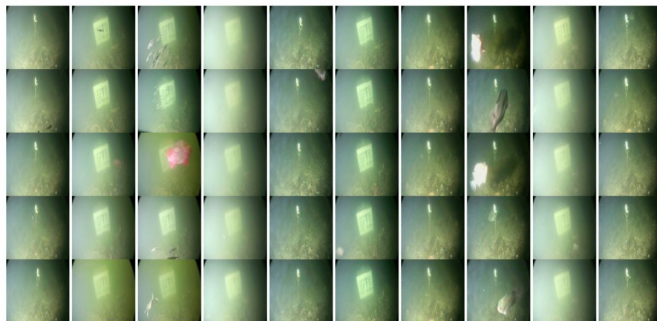
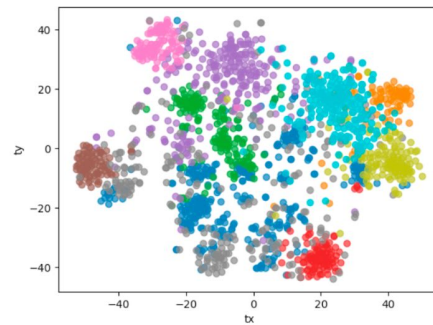
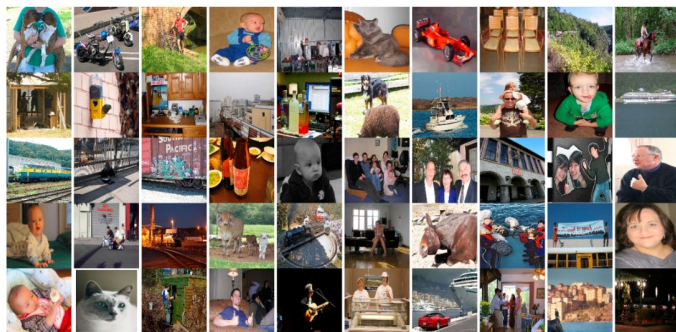
PVOC10



Brackish



# Clustering to get $H_2$



# Experiments and results

Table 6: YOLOv5s mean  $\bar{\theta}(\mathcal{M})$  and standard deviation  $\bar{\sigma}(\mathcal{M})$  on different PVOC10 partitioning.

| Partitioning | P (%)                   | R (%)                              | F1 (%)                  | mAP@.5 (%)              | mAP@.5:.95 (%)          |
|--------------|-------------------------|------------------------------------|-------------------------|-------------------------|-------------------------|
| KCV          | 61.33 $\pm$ 5.06        | 48.15 $\pm$ 3.05                   | 51.97 $\pm$ 3.04        | 50.13 $\pm$ 2.66        | 31.13 $\pm$ 2.72        |
| SKCV         | 59.62 $\pm$ 6.22        | 47.51 $\pm$ 3.47                   | 50.79 $\pm$ 3.19        | 49.52 $\pm$ 3.53        | 30.51 $\pm$ 3.00        |
| ISKCV        | 60.34 $\pm$ 5.02        | <b>48.20 <math>\pm</math> 2.35</b> | 52.02 $\pm$ 2.09        | 50.42 $\pm$ 2.65        | 31.31 $\pm$ 2.49        |
| FKCV         | 61.64 $\pm$ <b>4.68</b> | 47.95 $\pm$ 3.14                   | 51.64 $\pm$ <b>1.70</b> | 50.07 $\pm$ <b>2.24</b> | 31.27 $\pm$ <b>1.94</b> |
| Unified      | <b>62.63</b> $\pm$ 4.72 | 47.59 $\pm$ 2.77                   | <b>52.29</b> $\pm$ 2.31 | <b>50.43</b> $\pm$ 2.76 | <b>31.60</b> $\pm$ 2.91 |

Table 7: YOLOv5s mean  $\bar{\theta}(\mathcal{M})$  and standard deviation  $\bar{\sigma}(\mathcal{M})$  on different Brackish partitioning. \*Note that the biased split does not consider a video as an image subset and randomly splits all the images among the folds as shown in Figure 7.

| Partitioning  | P (%)                   | R (%)                          | F1 (%)                         | mAP@.5 (%)              | mAP@.5:.95(%)           |
|---------------|-------------------------|--------------------------------|--------------------------------|-------------------------|-------------------------|
| *Biased split | 98.10                   | 98.10                          | 98.10                          | 99.10                   | 83.60                   |
| KCV           | 63.24 $\pm$ 9.61        | 43.90 $\pm$ 6.75               | 48.64 $\pm$ 5.74               | 47.83 $\pm$ 7.16        | 29.39 $\pm$ 6.35        |
| SKCV          | 62.76 $\pm$ 9.58        | 44.55 $\pm$ 6.43               | 49.47 $\pm$ 6.26               | 48.97 $\pm$ 6.46        | 28.54 $\pm$ 3.62        |
| ISKCV         | 64.06 $\pm$ 10.01       | 46.02 $\pm$ 6.32               | 50.92 $\pm$ 7.19               | <b>50.25</b> $\pm$ 6.14 | <b>29.65</b> $\pm$ 3.73 |
| FKCV          | 62.94 $\pm$ <b>7.17</b> | 44.95 $\pm$ 5.33               | 49.15 $\pm$ 5.30               | 48.89 $\pm$ <b>5.54</b> | 29.43 $\pm$ 4.00        |
| Unified       | <b>66.38</b> $\pm$ 8.07 | <b>46.31</b> $\pm$ <b>5.22</b> | <b>51.53</b> $\pm$ <b>5.29</b> | 50.12 $\pm$ 5.98        | 29.54 $\pm$ <b>3.52</b> |

# Conclusion



- New paradigm for creating consistent K-fold in object detection tasks
- application to object detection
- results are significant
- 
- We adapted previous works on classification to object detection (S-KCV, D-SKCV, and DS-KCV) in a unified framework
- We showed interesting results both maximizing the performance estimation and minimizing the variance estimation, facilitating unbiased comparison of different algorithms on the same database