

# Fine-grained Recognition of Plants and Fungi from Images



Milan Šulc<sup>1</sup>



Jiří Matas<sup>1</sup>



Lukáš Pícek<sup>2</sup>

<sup>1</sup>Department of Cybernetics  
Faculty of Electrical Engineering  
Czech Technical University in Prague

<sup>2</sup>Department of Cybernetics  
Faculty of Applied Sciences  
University of West Bohemia

May 17, 2021 @ Computer Vision Methods

# Plant/Fungi Species Recognition



Many species identification tasks exist, depending on:

- ▶ the observation (view):
  1. from an image (or images) of a defined organ
  2. from image(s) of a larger part of the specimen
  3. from unspecified image(s) of the specimen
- ▶ acquisition conditions (controlled background, viewpoint, occlusion vs. unconstrained)
- ▶ number of *species* and granularity of the decision (typically fine-grained).

We formulate each of the considered tasks as a single-label classification problem on a set  $\mathcal{C} = \{1, 2, \dots, K\}$  of  $K$  classes (species). For fine-grained recognition  $K \gg 1$ .

# Fine-grained Species Recognition Is a Difficult Task

The observation depends on many factors:

- ▶ Genotype
- ▶ Age
- ▶ Season
- ▶ Local environment (Climate, Altitude, Illumination)
- ▶ Clutter (other plants in the foreground or background)
- ▶ Acquisition conditions
- ▶ Device

The classes (species) often have:

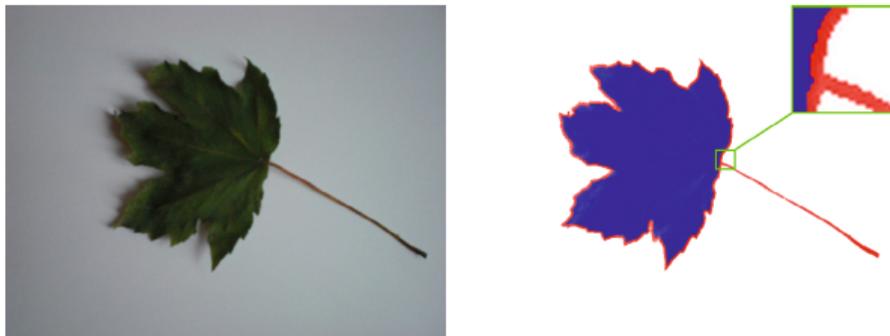
- ▶ Small inter-class differences
- ▶ High intra-class variability

# Texture Recognition Approach to Plant Recognition

Recognition of bark (bark texture covers all image area):



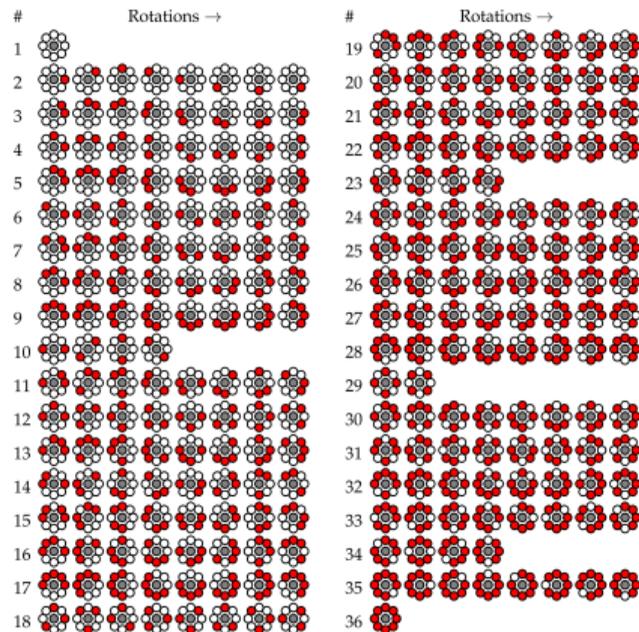
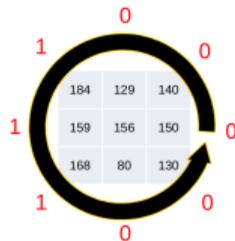
Recognition of segmented leaves:



# Fast Features Invariant to Rotation and Scale of Texture (Ffirst)

*Ffirst* extend previously published<sup>1,2</sup> texture descriptors.

- ▶ Completed Local Binary Patterns (sign- and magnitude-LBP)
- ▶ Rotation-invariant representation: “Histogram Fourier features”
- ▶ Scale space for multi-scale description and scale “invariance”
- ▶ Linear SVMs + explicit feature maps



<sup>1</sup>Milan Šulc and Jiří Matas. “Kernel-mapped histograms of multi-scale LBPs for tree bark recognition”. In: *IVCNZ 2013*.

<sup>2</sup>Milan Šulc and Jiří Matas. “Tree Identification from Images”. *Master Thesis*. Czech Technical University in Prague, 2014.

# Texture Recognition Approach: Results

- ▶ Best results on the Austrian Federal Forests bark dataset
- ▶ Best results on all 6 experimented leaf datasets, achieving  $>99\%$  accuracy on most datasets, including 99.5% accuracy on the MEW dataset of 153 tree species.<sup>3</sup>
- ▶ Competitive results<sup>4</sup> on standard texture recognition datasets:  
> 99% accuracy on the Brodatz32, UIUCTex, UMD, CURET and KTH-TIPS datasets, 87.9% and 76.6% acc. on the more difficult KTH-TIPS2a and KTI-TIPS2b datasets.<sup>5</sup>
- ▶ Fast descriptor: 200x200 px image takes about 0.05 sec. on a laptop (without a GPU).
- ▶ We noticed a significant color bias on several standard texture recognition datasets and proposed improvements to global color descriptors.

---

<sup>3</sup> Later outperformed by deep learning approach Šulc and Matas (Plant Methods, 2017), achieving  $> 99.9\%$  accuracy on the MEW dataset.

<sup>4</sup> Milan Šulc and Jiří Matas. "Fast Features Invariant to Rotation and Scale of Texture". In: *Computer Vision - ECCV 2014 Workshops (Zurich, Switzerland)*. Vol. 8926. LNCS. Springer, 2015, pp. 47–62. ISBN: 978-3-319-16180-8. DOI: 10.1007/978-3-319-16181-5\_4.

<sup>5</sup> The more recent CNN-based approach of Cimpoi et al. (IJCV, 2016) further improves the classification scores on most texture recognition datasets.

# Deep Learning Approach to Species Identification “in the Wild”

Deep Convolutional Neural Network (CNN) for the tasks of plant and fungi recognition “in the wild”:



- ▶ Large-scale databases available from computer vision challenges (LifeCLEF 2016-2019, FGVCx Flowers 2018, FGVCx Fungi 2018).
- ▶ **Prior shift** is a common phenomenon: Distribution of species on training data often differs from test data.

## Prior Shift

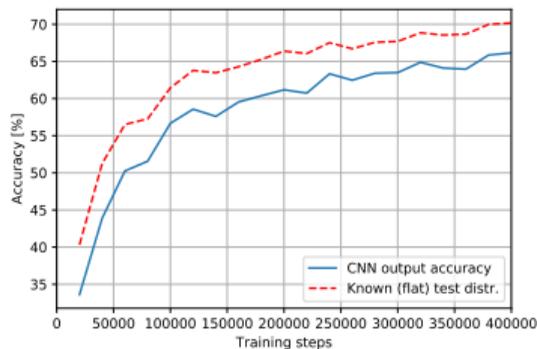
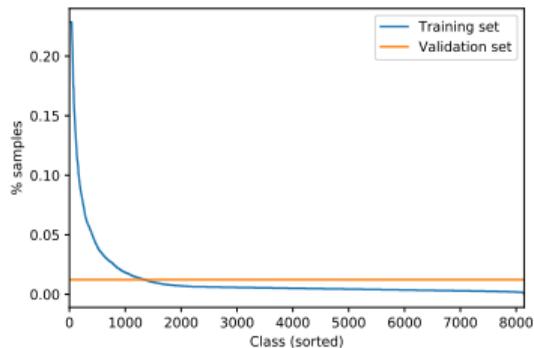
A CNN classifier  $f_{\text{CNN}} : \mathbb{R}^{W \times H \times 3} \mapsto \Delta^{K-1}$  is typically trained by minimizing the cross entropy loss  $L_{\text{CE}} = - \sum_{i=1}^N \log f_{\text{CNN}}(y_i | \mathbf{x}_i; \boldsymbol{\theta})$  on a training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  to estimate posterior probabilities:  $f_{\text{CNN}}(k | \mathbf{x}; \boldsymbol{\theta}) \approx p(k | \mathbf{x})$ .

Assume  $p(\mathbf{x} | k)$  remains unchanged, but class priors  $p_Y^e(k)$  on the new data differ from the training set  $p_Y(k)$ . Then  $p^e(k | \mathbf{x})$  differs from  $p(k | \mathbf{x}) \approx f_{\text{CNN}}(k | \mathbf{x})$ :

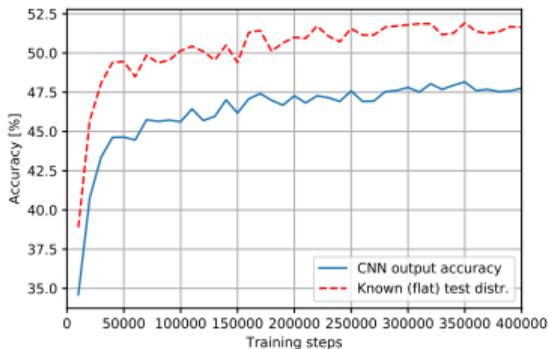
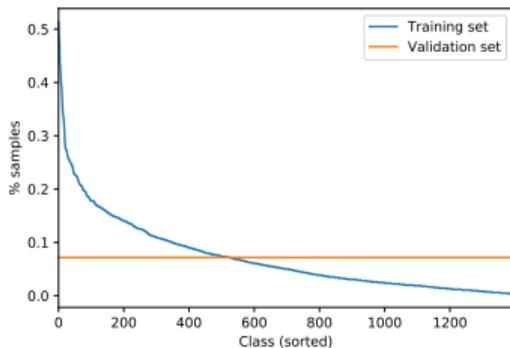
$$p^e(k | \mathbf{x}) = p(k | \mathbf{x}) \frac{p_Y^e(k) p_X(\mathbf{x})}{p_Y(k) p_X^e(\mathbf{x})} \propto p(k | \mathbf{x}) \frac{p_Y^e(k)}{p_Y(k)}$$

# When New Priors Are Known

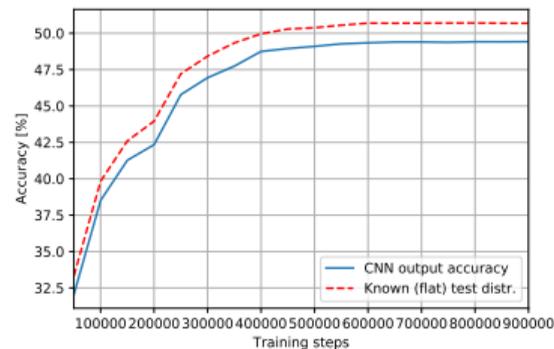
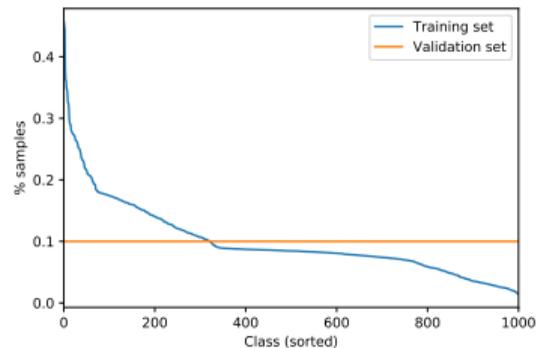
## iNaturalist 2018



## FGVCx Fungi 2018



## Webvision 1.0



## Maximum Likelihood Estimation of New Priors

Saerens et al.<sup>6</sup> proposed the following EM procedure to maximize the likelihood

$$L(\mathcal{E}) = \prod_{\mathbf{x} \in \mathcal{E}} p_X^e(\mathbf{x}) = \prod_{\mathbf{x} \in \mathcal{E}} \sum_{k=1}^K p^e(\mathbf{x}, k) = \prod_{\mathbf{x} \in \mathcal{E}} \sum_{k=1}^K p(\mathbf{x}|k) p_Y^e(k):$$

$$p^{(s)}(k|\mathbf{x}) = \frac{p(k|\mathbf{x}) \frac{\hat{p}_k^{(s)}}{p_Y(k)}}{\sum_{j=1}^K p(j|\mathbf{x}) \frac{\hat{p}_j^{(s)}}{p_Y(j)}} \quad (1)$$

$$\hat{p}_k^{(s+1)} = \frac{1}{N_e} \sum_{\mathbf{x} \in \mathcal{E}} p^{(s)}(k|\mathbf{x}) \quad (2)$$

---

<sup>6</sup>Marco Saerens et al. "Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure". In: *Neural computation* 14.1 (2002), pp. 21–41.

# Maximum Likelihood Estimation of New Priors

The EM algorithm may not result<sup>7</sup> in the unique optimal value.

We therefore also experiment with optimization of the MLE objective

$$\hat{\mathbf{p}}^{\text{MLE}} = \arg \max_{\mathbf{p}} \sum_{\mathbf{x} \in \mathcal{E}} \log \sum_{k=1}^K p_k \underbrace{\frac{p(k|\mathbf{x})}{p_Y(k)}}_{a_{ik}} p_X(\mathbf{x}) \quad \text{s.t.} \quad \sum_{k=1}^K p_k = 1; \quad \forall k : p_k \geq 0$$

by projected gradient ascent:

$$\hat{p}_k^{(s+1)} = \pi \left( \hat{p}_k^{(s)} + \lambda \frac{\partial \ell(\mathcal{E})}{\partial \hat{p}_k} \right), \quad \text{where} \quad \frac{\partial \ell(\mathcal{E})}{\partial \hat{p}_k} = \sum_{\mathbf{x} \in \mathcal{E}} \frac{a_{ik}}{\sum_{j=1}^K \hat{p}_j a_{ij}}$$

---

<sup>7</sup>Marthinus Christoffel Du Plessis and Masashi Sugiyama. "Semi-supervised learning of class balance under class-prior change by distribution matching". In: *Neural Networks* 50 (2014), pp. 110–119.

## Maximum A Posteriori Estimation of New Priors

Assuming a hyper-prior  $p(\mathbf{p})$  on the categorical priors:

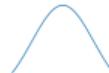
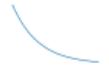
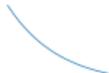
$$\begin{aligned}\hat{\mathbf{p}}^{\text{MAP}} &= \arg \max_{\mathbf{p}} p(\mathbf{p}|\mathcal{E}) = \arg \max_{\mathbf{p}} p(\mathbf{p}) \prod_{\mathbf{x} \in \mathcal{E}} p(\mathbf{x}|\mathbf{p}) = \\ &= \arg \max_{\mathbf{p}} \left[ \log p(\mathbf{p}) + \sum_{\mathbf{x} \in \mathcal{E}} \log p(\mathbf{x}|\mathbf{p}) \right] \text{ s.t. } \sum_{k=1}^K p_k = 1; \forall k : p_k \geq 0\end{aligned}$$

We use symmetric Dirichlet distribution:  $p(\mathbf{p}) = \frac{1}{B(\alpha)} \prod_{k=1}^K p_k^{\alpha-1}$   
favouring dense distributions, i.e. parameterized by  $\alpha > 1$ .

Log-concave, just adds gradient components:  $\frac{\partial \log p(\hat{\mathbf{p}})}{\partial \hat{p}_k} = \frac{\alpha - 1}{\hat{p}_k}$

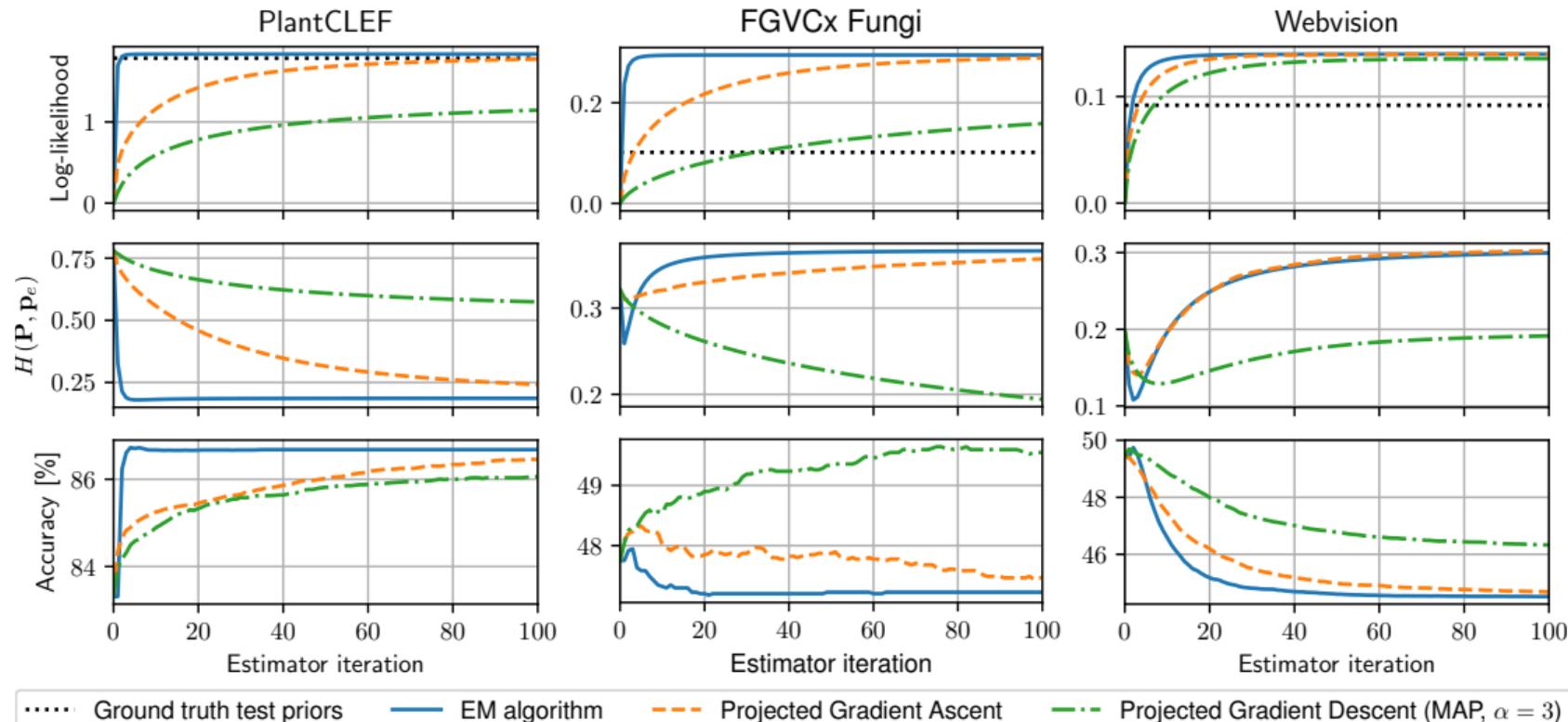
# Adjusting Predictions by Estimated Priors

CNN trained on unbalanced CIFAR-100 subsets and evaluated<sup>8</sup> on full test set:

Train. distribution												
Acc. (%)	48.15	55.70	60.88	64.01	65.62	<b>67.29</b>	36.68	47.72	54.00	56.57	60.37	61.66
after MLE	49.71	56.94	61.64	64.58	65.62	67.11	38.67	49.05	55.18	57.05	60.59	61.74
after MAP, $\alpha = 3$	49.75	56.94	61.65	<b>64.59</b>	65.64	67.18	38.75	49.20	55.19	<b>57.10</b>	60.58	<b>61.76</b>
after MAP, $\alpha = 10$	<b>50.07</b>	<b>56.97</b>	<b>61.68</b>	64.55	<b>65.70</b>	67.23	<b>39.12</b>	<b>49.34</b>	<b>55.22</b>	<b>57.10</b>	<b>60.69</b>	<b>61.76</b>
with known $p_Y^e(k)$	51.20	57.61	62.23	64.73	65.92	67.44	40.62	50.07	55.86	57.49	60.92	62.11

<sup>8</sup>Milan Šulc and Jiří Matas. "Improving CNN classifiers by estimating test-time priors". In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2019.

# Adjusting Predictions by Estimated Priors



# LifeCLEF Plant (Herbs, Trees, Ferns) Identification Challenges

## PlantCLEF 2016 →→

- ▶ 1 000 plant species, 113K training images

## PlantCLEF 2017

- ▶ 10,000 plant species
- ▶ “trusted” train. set: 256K images (EoL.org)
- ▶ “noisy” train. set: 1.4M images (web search)

## ExpertLifeCLEF 2018

- ▶ same 10,000 species as in 2017
- ▶ comparison against human experts

## PlantCLEF 2019

- ▶ 10,000 species from “data-deficient” regions
- ▶ “trusted” train. set: 72K images (EoL.org)
- ▶ “noisy” train. set: 375K images (web search)



## Meta-data available:

- ▶ “types of view”: *Leaf, Leaf Scan, Flower, Fruit, Stem, Branch, Entire*
- ▶ “Observation ID” available

# LifeCLEF Plant (Herbs, Trees, Ferns) Identification Challenges

## PlantCLEF 2016

- ▶ 1 000 plant species, 113K training images

## PlantCLEF 2017 →→

- ▶ 10,000 plant species
- ▶ “trusted” train. set: 256K images (EoL.org)
- ▶ “noisy” train. set: 1.4M images (web search)

## ExpertLifeCLEF 2018 →→

- ▶ same 10,000 species as in 2017
- ▶ comparison against human experts

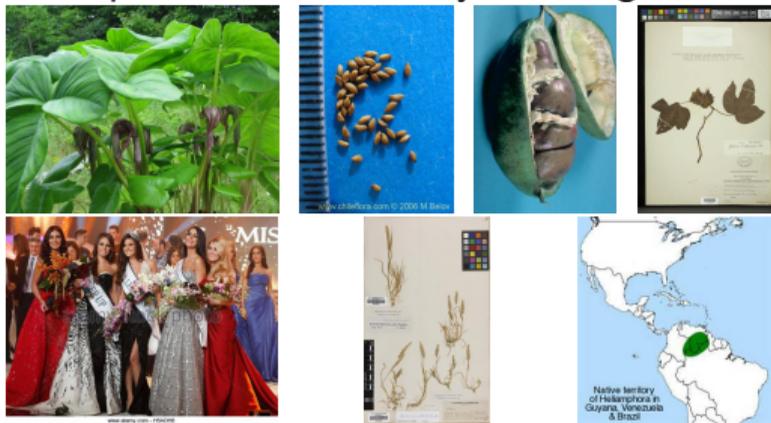
## PlantCLEF 2019

- ▶ 10,000 species from “data-deficient” regions
- ▶ “trusted” train. set: 72K images (EoL.org)
- ▶ “noisy” train. set: 375K images (web search)

## Examples from the “trusted” training set:



## Examples from the “noisy” training set:



# LifeCLEF Plant (Herbs, Trees, Ferns) Identification Challenges

## PlantCLEF 2016

- ▶ 1 000 plant species, 113K training images

## PlantCLEF 2017

- ▶ 10,000 plant species
- ▶ “trusted” train. set: 256K images (EoL.org)
- ▶ “noisy” train. set: 1.4M images (web search)

## ExpertLifeCLEF 2018

- ▶ same 10,000 species as in 2017
- ▶ comparison against human experts

## PlantCLEF 2019 →→

- ▶ 10,000 species from “data-deficient” regions
- ▶ “trusted” train. set: 72K images (EoL.org)
- ▶ “noisy” train. set: 375K images (web search)

- ▶ Species from the Guiana shield and the Amazon rain forest
- ▶ 20% of species: less than 10 images

Examples from the “trusted” training set:



Examples of wrongly annotated images from the “noisy” training set:



# FGVC 2018 Competitions

## FGVCx Fungi:

- ▶ 1394 fungi species
- ▶ 85 578 train. + 4 182 val. images
- ▶ 9 758 competition test images



## FGVCx Flower classification challenge:

- ▶ 997 species of flowering plants
- ▶ 669 304 train. images, no val. set
- ▶ 12 961 competition test images



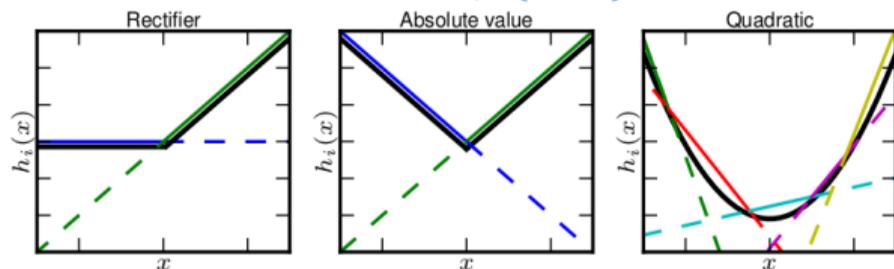
## iNaturalist 2018:

- ▶ 8 142 species (plants, fungi, animals, ...)
- ▶ 437 513 train. + 24 426 val. images
- ▶ 149 394 competition test images



# Species Recognition Competitions: Methods

- ▶ Fine-tuning state-of-the-art CNN classifiers<sup>9,10</sup> from ImageNet checkpoints
- ▶ Combining predictions from an ensemble of CNNs (sum or mode)
- ▶ Combining predictions from images with the same ObservationID (if available)
- ▶ Filtering noisy data
- ▶ Image augmentation (during training and at test time)
- ▶ Running averages (exponential decay) of trainable variables:  
68.4%  $\rightarrow$  80.2% acc. (!) with Inception-ResNet-v2
- ▶ Adjusting predictions to prior shift
- ▶ Maxout<sup>11</sup>:  $\forall i \in \{1, \dots, m\} : h_i(x) = \max_{j \in \{1, \dots, s\}} z_{ij}, z_{ij} = \mathbf{x}^\top \mathbf{W}_{.ij} + b_{ij}$



<sup>9</sup>Kaiming He et al. "Deep residual learning for image recognition". In: *CVPR*. 2016.

<sup>10</sup>Christian Szegedy et al. "Inception-v4, inception-resnet and the impact of residual connections on learning.". In: *AAAI*. 2017.

# Species Recognition Competitions: What Didn't Work

- ▶ Handling "types of view" separately decreased precision in preliminary experiments
- ▶ "Bootstrapping"<sup>12</sup> consistency objectives for training on noisy labels

$$L_{\text{soft}}(\mathbf{q}, \mathbf{t}) = \sum_{k=1}^K [\beta t_k + (1 - \beta) q_k] \log q_k, \quad L_{\text{hard}}(\mathbf{q}, \mathbf{t}) = \sum_{k=1}^K [\beta t_k + (1 - \beta) z_k] \log q_k$$

- ▶ *soft bootstrapping* uses the softmax predictions  $q_k$ ,
- ▶ *hard bootstrapping* uses only the strongest prediction  $z_k = \begin{cases} 1 & \text{if } k = \arg \max_i q_i \\ 0 & \text{otherwise} \end{cases}$

---

<sup>12</sup>Scott Reed et al. "Training deep neural networks on noisy labels with bootstrapping". In: *arXiv preprint arXiv:1412.6596* (2014).

# LifeCLEF Competitions: Results

	Place	Our score	Winner score	Metric
PlantCLEF 2016	3rd	70.1% mAP	74.2% mAP	Mean average precision
+ ObservationID $\Sigma$		78.8% mAP		
PlantCLEF 2017	3rd	84.3% MRR	92.0% MRR	Mean reciprocal rank
+ correct prior shift		86.7% MRR		$\text{MRR} = \frac{1}{ Q } \sum_{i=1}^{ Q } \frac{1}{\text{rank}_i}$
PlantCLEF 2018	<b>1st</b>	88.4% acc.	=	Top-1 accuracy
PlantCLEF 2019	(post-challenge)	31.9% acc.	24.7%	Top-1 accuracy

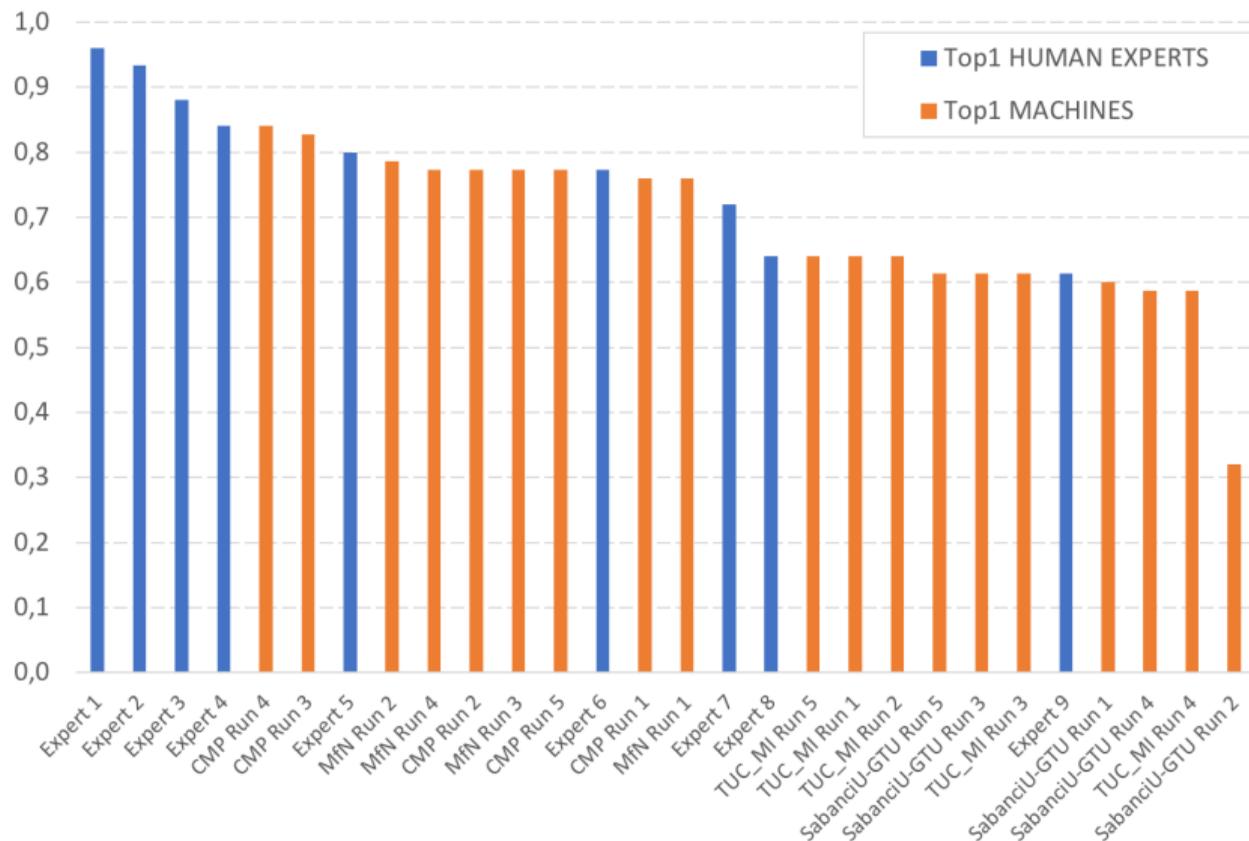
# FGVC 2018 Competitions: Official Results (top 10)

FGVCx Fungi		Recall@3 Err.(%)	
Team Name	Private	Public	
<b>CMP (ours)</b>	<b>21.197</b>	<b>20.772</b>	
digitalspecialists	23.188	23.471	
Val An	25.091	25.213	
DL Analytics	28.341	26.853	
Invincibles	28.751	28.493	
Tian Xi	32.235	31.636	
Igor Krashenyi	32.616	34.164	
wakaka	42.219	41.339	
George Yu	47.621	47.113	
Xinshao	67.837	67.509	

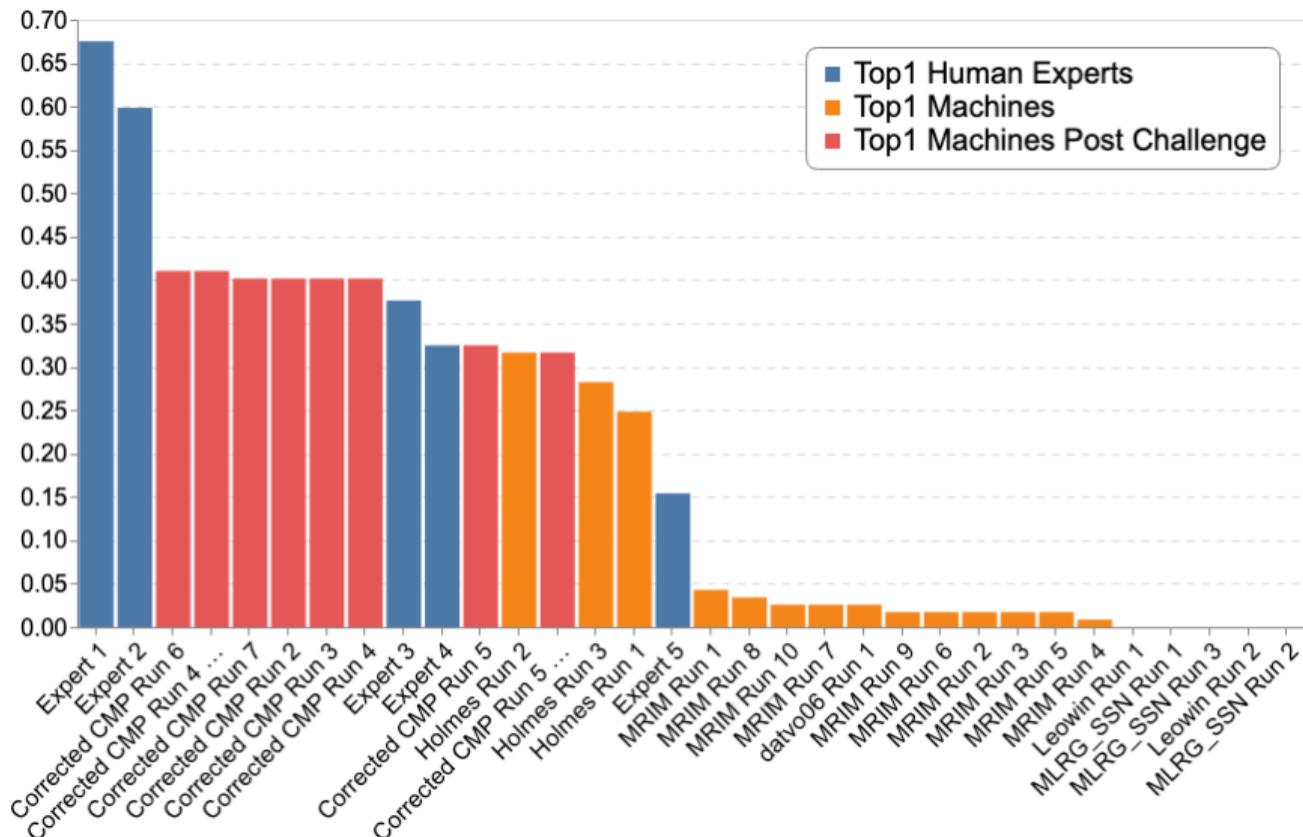
iNaturalist		Recall@3 Err.(%)	
Team Name	Private	Public	
DLUT VLG	12.858	13.068	
DL Analytics	13.981	14.214	
fadivugibs	14.618	14.914	
<b>CMP (ours)</b>	<b>16.076</b>	<b>16.360</b>	
fISHpAM	16.892	17.149	
traveler	16.988	17.235	
yen	17.201	17.412	
Shuang	18.357	18.549	
Mr.M	20.092	20.291	
Dequan Wang	20.814	21.157	

FGVCx Flowers		Top 1 Err.(%)	
Team Name	Private	Public	
<b>CMP (ours)</b>	<b>7.599</b>	<b>6.828</b>	
fadivugibs	8.177	7.638	
DLUT VLG	8.242	7.677	
yen	8.396	7.716	
xiaoxiao	9.579	8.641	
NDer MJU	11.636	10.976	
thesouthfrog	15.211	13.618	
nimahai	16.368	15.258	
Miroslav Štola	20.187	19.637	
Imao	20.342	20.177	

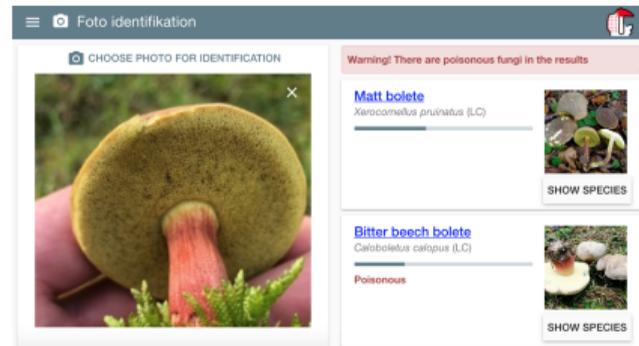
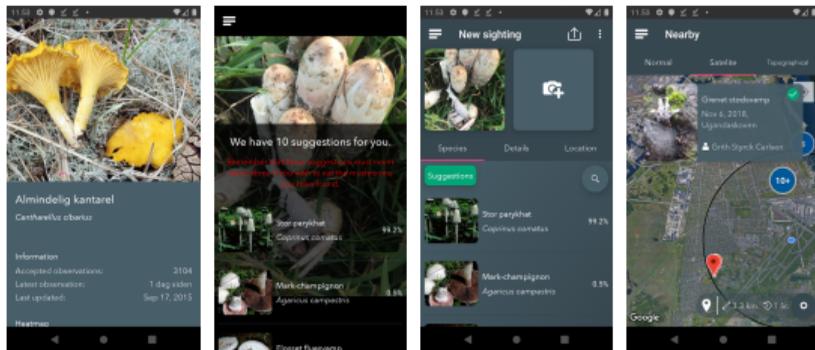
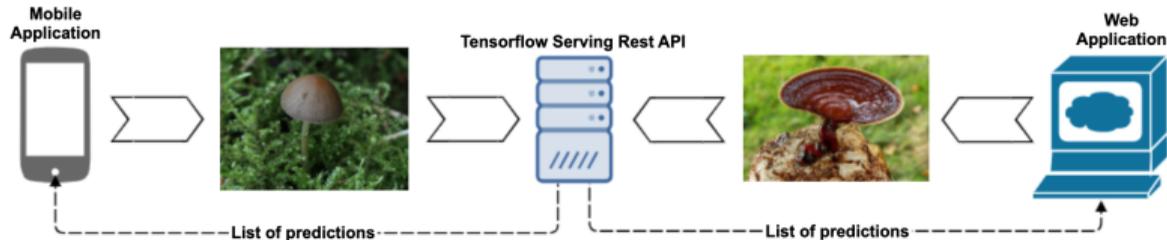
# ExpertLifeCLEF 2018: Experts vs. Machines



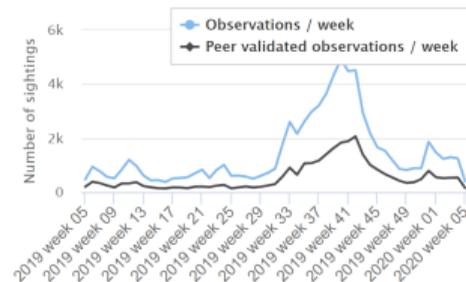
# PlantCLEF 2019: Experts vs. Machines on Data Deficient Species



# Fungi Recognition Service for the Atlas of Danish Fungi



Oct.-Dec. 2018: 11K new sightings with images  
Oct.-Dec. 2019: 21K new sightings with images  
Oct.-Dec. 2020: 28K new sightings with images



# Danish Fungi 2020 Dataset

- ▶ fungi observations from 1874 till the end of 2020
- ▶ precise annotation with taxonomy-accurate class labels
- ▶ highly unbalanced long-tailed class distribution
- ▶ rich observation meta-data (time, location, substrate, habitat, .. for  $> 99\%$  images)
- ▶ zero overlap with ImageNet

## **DF20**

- ▶ 1,604 classes
- ▶ 248,466 training images
- ▶ 27,608 (public) test images

## **DF20 – Mini**

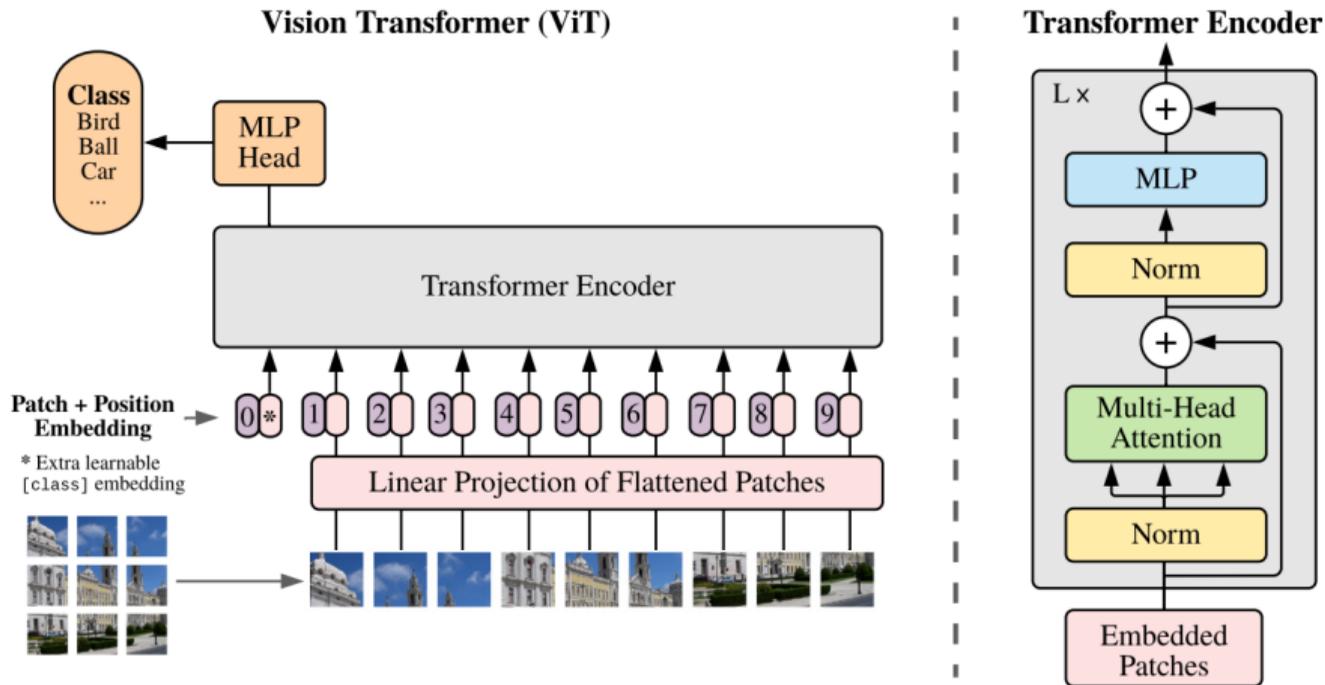
- ▶ 182 classes
- ▶ 32,753 training images
- ▶ 3,640 (public) test images

# Danish Fungi 2020: Baseline Experiments

			Top1 [%]	Top3 [%]	$F_1^m$			Top1 [%]	Top3 [%]	$F_1^m$
DF20-Mini	EfficientNet-B0	$224 \times 224$	65.66	83.35	0.531	$384 \times 384$	70.22	85.69	0.596	
	EfficientNet-B3		66.90	83.49	0.537		72.09	87.17	0.624	
	SE-ResNeXt-101		69.48	85.58	0.593		72.34	87.53	0.631	
	ViT-Base/16		69.37	<b>86.54</b>	0.589		74.84	88.74	0.655	
	ViT-Large/16		<b>70.71</b>	86.51	<b>0.599</b>		<b>75.96</b>	<b>89.37</b>	<b>0.664</b>	
DF20	EfficientNet-B0	$224 \times 224$	70.38	85.18	0.613	$384 \times 384$	75.27	88.65	0.670	
	ViT-Base/16		73.45	87.15	0.658		79.40	90.93	0.724	
	ViT-Large/16		<b>75.32</b>	<b>88.12</b>	<b>0.679</b>		<b>81.25</b>	<b>91.93</b>	<b>0.747</b>	

**Table 1:** Classification results of selected CNN and ViT architectures on DF20 and DF20-Mini dataset for two input resolutions.

# Best on Danish Fungi 2020: Vision Transformers<sup>13</sup>



<sup>13</sup>Alexey Dosovitskiy et al. "An image is worth 16x16 words: Transformers for image recognition at scale". In: *arXiv preprint arXiv:2010.11929* (2020).

## Danish Fungi 2020: Baseline for Utilizing Metadata

Assumption:  $P(I|S) = P(I|S, D)$ .

Then the class posterior given the image  $I$  and metadata  $D$  is:

$$P(S|I, D) = P(S|I) \frac{P(S|D)}{P(S)} \frac{P(I)}{P(I|D)} \propto P(S|I) \frac{p(S|D)}{p(S)},$$

To use several metadata, we assume:

$$P(S|D_1, D_2) \propto \frac{P(S|D_1)P(S|D_2)}{P(S)}.$$

# Danish Fungi 2020: Utilizing Metadata

	<b>H</b>	<b>M</b>	<b>S</b>	<b>Top1 [%]</b>	<b>Top3 [%]</b>	<b>F<sub>1</sub><sup>m</sup></b>
Danish Fungi 2020	×	×	×	73.45	87.15	0.658
	✓	×	×	+2.00	+1.42	+0.036
	×	✓	×	+1.37	+1.23	+0.024
	×	×	✓	+0.98	+0.96	+0.016
	×	✓	✓	+2.30	+2.10	+0.039
	✓	×	✓	+2.92	+2.41	+0.051
	✓	✓	×	+3.16	+2.50	+0.056
	✓	✓	✓	<b>+3.58</b>	<b>+3.05</b>	<b>+0.062</b>

**Table 2:** Performance gains from *Fungus* observation metadata: H - Habitat, S - Substrate, M - Month, and their combinations, on DF20. ViT-Base/16 with image size  $224 \times 224$ .

## The Danish Fungi 2020 dataset and benchmark

- ▶ published on arXiv<sup>14</sup>
- ▶ submitted to a major computer vision conference (in review)
- ▶ coming soon: public benchmark with **DNA sequenced private test set**

## Classifier calibration and prior shift adaptation

- ▶ recently submitted a patent application with Toyota
- ▶ submitted to a major computer vision conference (in review)

---

<sup>14</sup>Lukáš Pícek et al. "Danish Fungi 2020 – Not Just Another Image Recognition Dataset". In: *arXiv preprint arXiv:2103.10107* (2021).

## Estimating new priors using Confusion Matrix

A standard procedure [13,15] for prior estimation is based on the classifier confusion matrix in the format  $C_{d|y}$ , containing probabilities  $p(D = i|Y = k)$  of deciding for label  $i$  when the true class is  $k$ .

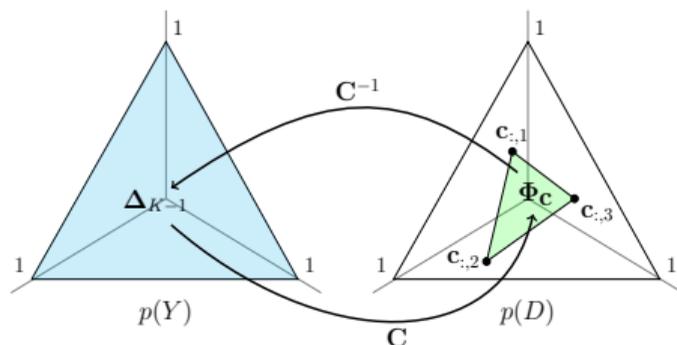
$$p(D = i) = \sum_{k=1}^K p(D = i|Y = k)p(Y = k)$$
$$p(D) = C_{d|y}p(Y)$$

The new priors can then be estimated as follows:

$$\hat{p}_{\mathcal{E}}(Y) = \hat{C}_{d|y}^{-1}\hat{p}_{\mathcal{E}}(D)$$

## Problem of the simple Confusion Matrix method

The confusion matrix defines a convex set  $\Phi_C$  of feasible values  $p(D)$  within the probability simplex:



Inconsistent estimates of the Confusion Matrix and  $p(D)$  can result in a vector outside of the probability simplex, i.e. the estimate can contain negative values. E.g.:

$$\hat{C}_{d|y} = \begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{bmatrix}, \hat{p}_{\mathcal{E}}(D) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \rightarrow \hat{p}(Y) = \begin{bmatrix} 4/3 \\ -1/3 \end{bmatrix}$$

## The proposed MLE based on Confusion Matrix

We assume a multinomial distribution of the observed classifier decision:

$$L(\mathbf{Q}) = p(\mathbf{n}|\mathbf{Q}) = \frac{(n_1 + \dots + n_K)!}{n_1! \cdot \dots \cdot n_K!} \cdot q_1^{n_1} \cdot \dots \cdot q_K^{n_K}$$

The log-likelihood  $\ell(\mathbf{P}) = \log p(\mathbf{n}|\mathbf{P}) = \sum_{k=1}^K n_k \log(\mathbf{c}_{k,:} \cdot \mathbf{P}) + \theta_{\mathbf{n}}$  is then maximized with condition on solutions on the probabilistic simplex:

$$\begin{aligned} \hat{\mathbf{P}} &= \arg \max_{\mathbf{P}} \ell(\mathbf{P}) = \arg \max_{\mathbf{P}} \sum_{k=1}^K n_k \log \mathbf{c}_{k,:} \cdot \mathbf{P} \\ \text{s.t.} &: \sum_{k=1}^K P_k = 1; \quad \forall k : P_k \geq 0 \end{aligned}$$

## Adding a hyper-prior is convenient

$$\begin{aligned}\hat{P}_{MAP} &= \arg \max_{\mathbf{P}} p(\mathbf{P}|\mathbf{n}) = \arg \max_{\mathbf{P}} p(\mathbf{P})p(\mathbf{n}|\mathbf{P}) \\ &= \arg \max_{\mathbf{P}} \log p(\mathbf{P}) + \arg \max_{\mathbf{P}} \log p(\mathbf{n}|\mathbf{P}) \\ \text{s.t.} \quad &\forall k : P_k \geq 0; \quad \sum_{k=1}^K P_k = 1\end{aligned}$$

# Results

Dataset		BCTS calibrated	NA	CM	CM <sup>L</sup>	SCM	SCM <sup>L</sup>	Oracle
CIFAR100*	LT→UNI	✗	31.66 <sup>±1.27</sup>	22.02 <sup>±2.56</sup>	<b>33.00<sup>±1.56</sup></b>	26.64 <sup>±3.85</sup>	<b>33.50<sup>±1.38</sup></b>	34.06 <sup>±1.35</sup>
	LT→UNI	✓	31.71 <sup>±1.29</sup>	18.56 <sup>±3.70</sup>	<b>31.90<sup>±1.49</sup></b>	24.52 <sup>±5.49</sup>	<b>33.78<sup>±1.30</sup></b>	34.54 <sup>±1.32</sup>
	UNI→LT	✗	63.83 <sup>±0.82</sup>	<b>68.06<sup>±0.92</sup></b>	68.04 <sup>±0.76</sup>	68.10 <sup>±0.81</sup>	<b>68.23<sup>±0.71</sup></b>	69.13 <sup>±0.58</sup>
	UNI→LT	✓	63.83 <sup>±0.82</sup>	<b>69.08<sup>±0.94</sup></b>	69.07 <sup>±1.00</sup>	69.31 <sup>±0.94</sup>	<b>69.38<sup>±0.78</sup></b>	70.65 <sup>±0.75</sup>
Places365	LT→UNI	✗	25.14 <sup>±0.14</sup>	17.79 <sup>±0.28</sup>	<b>27.77<sup>±0.45</sup></b>	19.78 <sup>±2.21</sup>	<b>28.64<sup>±0.22</sup></b>	32.99 <sup>±0.46</sup>
	LT→UNI	✓	25.16 <sup>±0.14</sup>	15.52 <sup>±0.66</sup>	<b>27.11<sup>±0.32</sup></b>	17.86 <sup>±1.62</sup>	<b>27.32<sup>±0.18</sup></b>	33.51 <sup>±0.25</sup>
	UNI→LT	✗	58.17 <sup>±1.01</sup>	81.16 <sup>±0.61</sup>	<b>81.67<sup>±0.66</sup></b>	82.04 <sup>±0.15</sup>	<b>82.09<sup>±0.64</sup></b>	88.14 <sup>±0.27</sup>
	UNI→LT	✓	58.17 <sup>±1.01</sup>	81.16 <sup>±0.59</sup>	<b>81.70<sup>±0.61</sup></b>	82.04 <sup>±0.15</sup>	<b>82.42<sup>±0.90</sup></b>	88.15 <sup>±0.30</sup>
ImageNet	LT→UNI	✗	34.30 <sup>±0.19</sup>	19.48 <sup>±0.74</sup>	<b>33.57<sup>±0.33</sup></b>	23.94 <sup>±2.04</sup>	<b>35.78<sup>±0.17</sup></b>	37.34 <sup>±0.15</sup>
	LT→UNI	✓	34.31 <sup>±0.19</sup>	16.78 <sup>±0.64</sup>	<b>31.94<sup>±0.51</sup></b>	20.26 <sup>±4.08</sup>	<b>35.55<sup>±0.23</sup></b>	37.45 <sup>±0.19</sup>

**Table 3:** “Improve Estimates from Confusion Matrix.” Accuracy ( $\pm$  std. dev.) after adaptation with new prior estimate based on confusion matrix (CM) inversion and our proposed method (CM<sup>L</sup>). SCM denotes soft confusion matrix, NA denotes no adaptation, Oracle is adaptation with ground truth priors. Results on CIFAR are averaged from 10 experiments, results on Places and ImageNet are averaged from 5 experiments.

# Results

Dataset		BCTS calibrated	NA	EM	MLE CM <sup>L</sup>	SCM <sup>L</sup>	MAP	MAP CM <sup>M</sup>	SCM <sup>M</sup>	Oracle
CIFAR100*	LT→UNI	✗	31.66 $\pm$ 1.27	32.81 $\pm$ 1.41	33.00 $\pm$ 1.56	<b>33.50<math>\pm</math>1.38</b>	32.87 $\pm$ 1.40	33.88 $\pm$ 1.40	<b>33.92<math>\pm</math>1.36</b>	34.06 $\pm$ 1.35
	LT→UNI	✓	31.71 $\pm$ 1.29	26.18 $\pm$ 1.61	31.90 $\pm$ 1.49	<b>33.78<math>\pm</math>1.30</b>	27.85 $\pm$ 1.54	33.96 $\pm$ 1.46	<b>34.11<math>\pm</math>1.62</b>	34.54 $\pm$ 1.32
	UNI→LT	✗	63.83 $\pm$ 0.82	67.23 $\pm$ 0.88	68.04 $\pm$ 0.76	<b>68.23<math>\pm</math>0.71</b>	66.89 $\pm$ 0.92	67.11 $\pm$ 0.91	<b>67.13<math>\pm</math>0.79</b>	69.13 $\pm$ 0.58
	UNI→LT	✓	63.83 $\pm$ 0.82	69.17 $\pm$ 0.91	69.07 $\pm$ 1.00	<b>69.38<math>\pm</math>0.78</b>	68.30 $\pm$ 0.77	<b>68.42<math>\pm</math>0.85</b>	68.26 $\pm$ 0.72	70.65 $\pm$ 0.75
Places365	LT→UNI	✗	25.14 $\pm$ 0.14	28.02 $\pm$ 0.92	27.77 $\pm$ 0.45	<b>28.64<math>\pm</math>0.22</b>	28.86 $\pm$ 0.68	30.27 $\pm$ 0.41	<b>30.88<math>\pm</math>0.36</b>	32.99 $\pm$ 0.46
	LT→UNI	✓	25.16 $\pm$ 0.14	26.25 $\pm$ 1.02	27.11 $\pm$ 0.32	<b>27.32<math>\pm</math>0.18</b>	28.29 $\pm$ 0.08	29.37 $\pm$ 0.30	<b>29.39<math>\pm</math>0.14</b>	33.51 $\pm$ 0.25
	UNI→LT	✗	58.17 $\pm$ 1.01	<b>82.63<math>\pm</math>0.31</b>	81.67 $\pm$ 0.66	82.09 $\pm$ 0.64	<b>78.81<math>\pm</math>0.47</b>	76.97 $\pm$ 0.45	75.54 $\pm$ 0.43	88.14 $\pm$ 0.27
	UNI→LT	✓	58.17 $\pm$ 1.01	<b>82.70<math>\pm</math>0.34</b>	81.70 $\pm$ 0.61	82.42 $\pm$ 0.90	<b>77.05<math>\pm</math>0.44</b>	76.12 $\pm$ 0.59	73.29 $\pm$ 0.46	88.15 $\pm$ 0.30
ImageNet	LT→UNI	✗	34.30 $\pm$ 0.19	34.63 $\pm$ 0.29	33.57 $\pm$ 0.33	<b>35.78<math>\pm</math>0.17</b>	35.10 $\pm$ 0.21	36.51 $\pm$ 0.17	<b>36.90<math>\pm</math>0.16</b>	37.34 $\pm$ 0.15
	LT→UNI	✓	34.31 $\pm$ 0.19	24.39 $\pm$ 2.31	31.94 $\pm$ 0.51	<b>35.55<math>\pm</math>0.23</b>	32.08 $\pm$ 2.11	35.39 $\pm$ 0.09	<b>36.80<math>\pm</math>0.16</b>	37.45 $\pm$ 0.19

**Table 4:** “How to estimate new priors?” Accuracy ( $\pm$  std. dev.) after adaptation to new priors estimated with different Maximum Likelihood and Maximum A Posteriori estimates. NA denotes no adaptation, Oracle is adaptation with ground truth priors. Results on CIFAR are averaged from 10 experiments, results on Places and ImageNet are averaged from 5 experiments.

# Summary

- ▶ *Fast Features Invariant to Rotation and Scale of Texture*
  - ▶ accuracy  $> 99\%$  on most leaf datasets
  - ▶ competitive results in texture classification:  $> 99\%$  acc. on standard texture datasets
  - ▶ computationally efficient: 200x200 px image takes about 0.05s
- ▶ Deep learning approach for “in-the-wild” species classification:
  - ▶ competition-winning contributions to fine-grained recognition challenges
  - ▶ reached human-expert level of accuracy
  - ▶ applied as an online recognition service for a Atlas of Danish Fungi, increasing the involvement of users in biodiversity data collection
- ▶ Adjusting to prior shift is important.
  - ▶ estimating new priors using MLE compared with a proposed MAP estimation, increasing the reliability and accuracy in several tasks
  - ▶ an analogous approach can be used to utilize metadata in the decision process
- ▶ The Danish Fungi 2020 dataset for fine-grained classification:
  - ▶ accurate class labels, zero overlap with ImageNet, and rich meta-data

Thank you.

Discussion.