# Methodology of automatic classification of plant pollen using machine learning

**Artur Tomczak[1], J. Vasilescu[2], I.S. Stachlewska[1], Ł. Janicka[1]**

[1] Faculty of Physics, University of Warsaw (UW), 02-093 Warsaw, Poland
[2] National Institute of Research and Development for Optoelectronics INOE 2000, 077125 Magurele, Romania

*Contact emails:*        *artur.tomczak@fuw.edu.pl*
                          *iwona.stachlewska@fuw.edu.pl*

Remote Sensing Laboratory

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

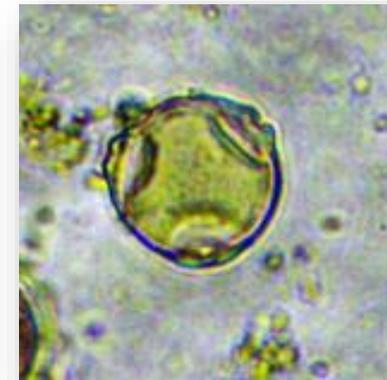**Motivation**: *Health impact of allergenic pollen.*

**Objective**: *Investigate allergenic pollen in a real time over an urban continental site.*

**Innovation**: *open-source (free) and flexible algorithm that allows to investigate, train and predict pollen taxa.*
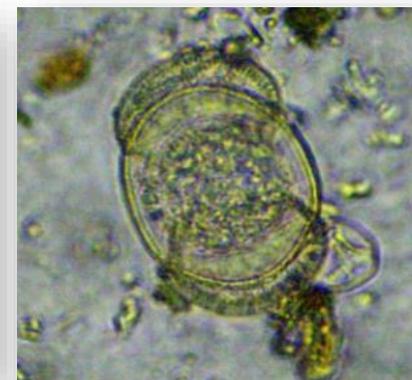


Sensor Rapid-E

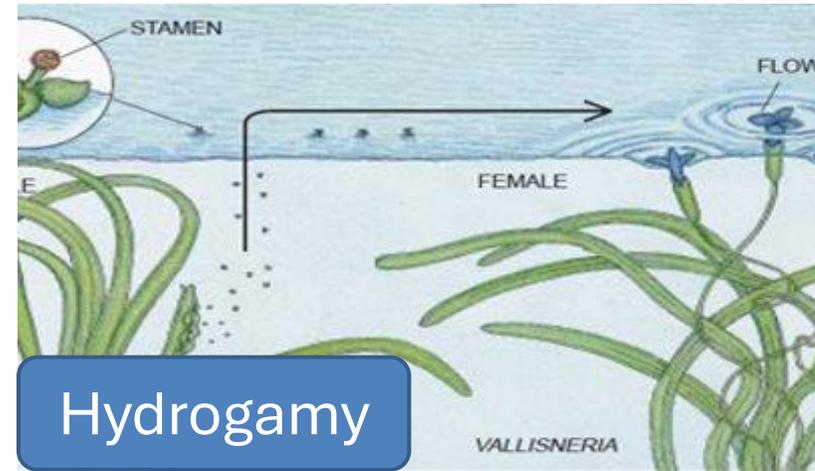*Betula* pollen grain     Pine pollen grain

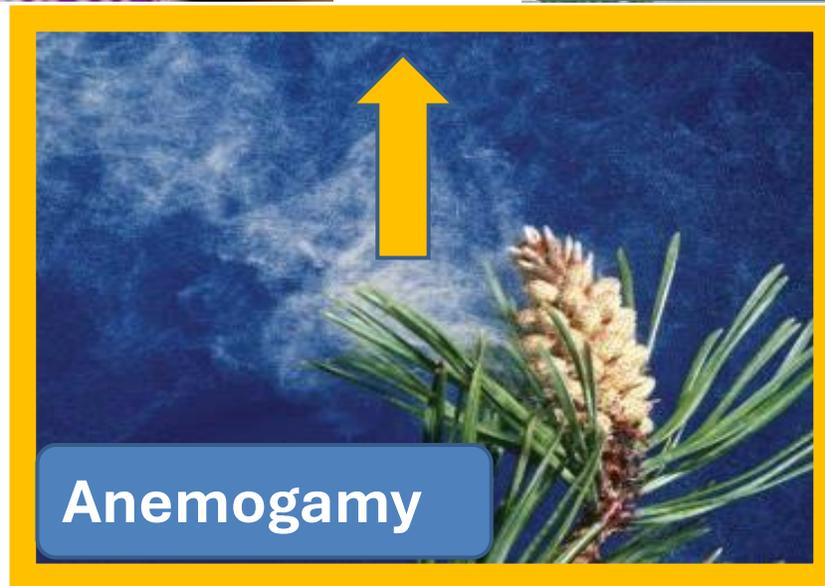*Photo provided by Zuzanna Rykowska*
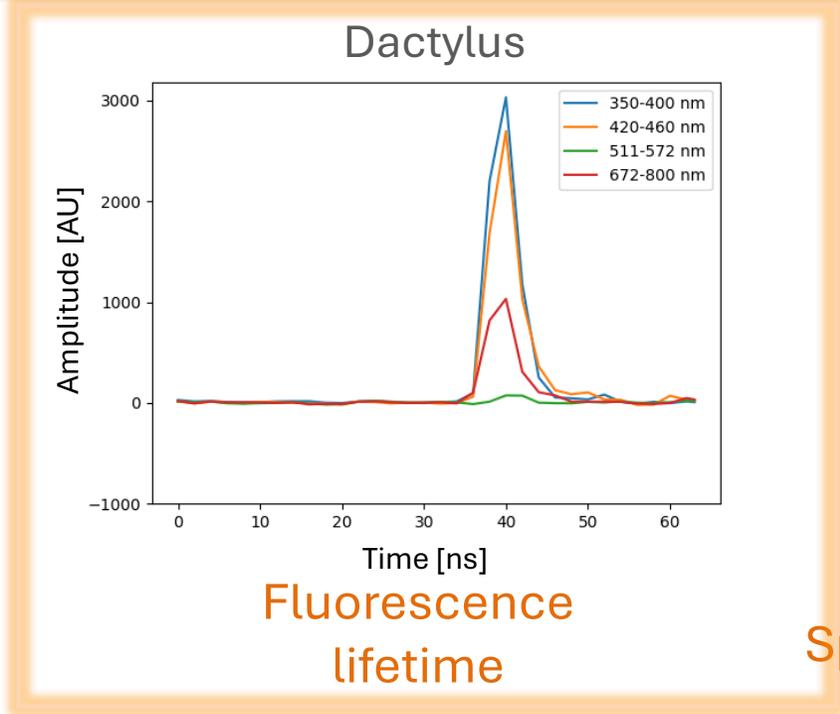
Zoogamy

Hydrogamy

Anemogamy

- Atmospheric dynamics
- Optical properties
- Microphysical properties

# Rapid-E sampler

- Ranges:
  - 350-400 nm
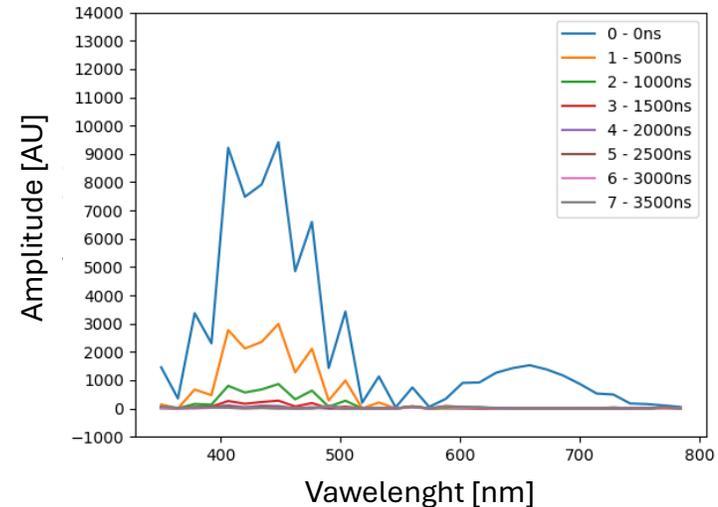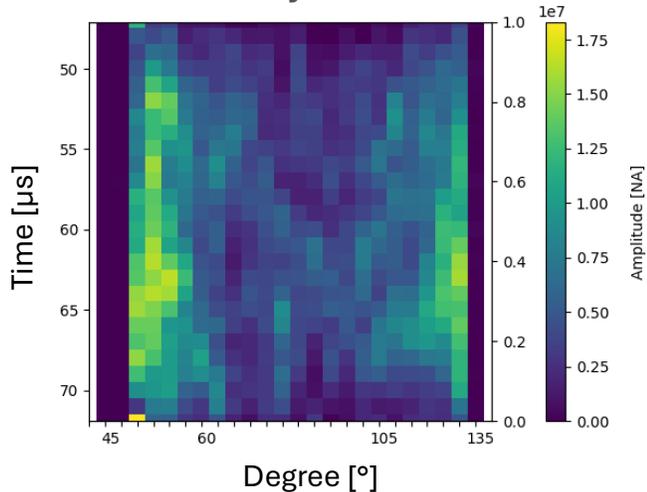  - 420-460 nm
  - 511-572 nm
  - 672-800 nm


Dactylus

Sensor Rapid-E

**Scattering**

**Fluorescence lifetime**

**Spectral ranges**
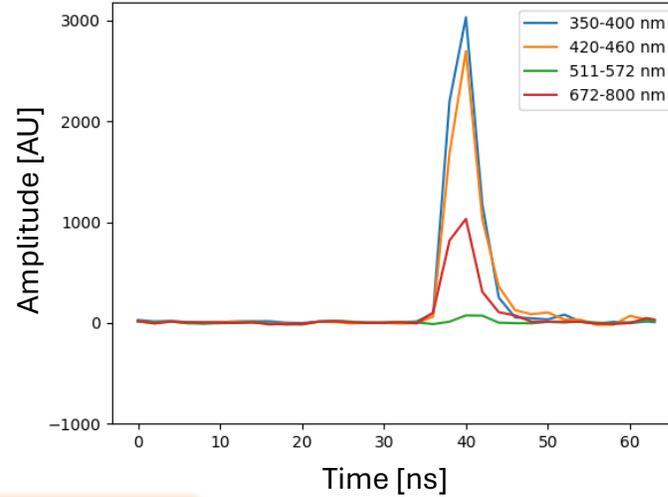

Corylus


Alnus

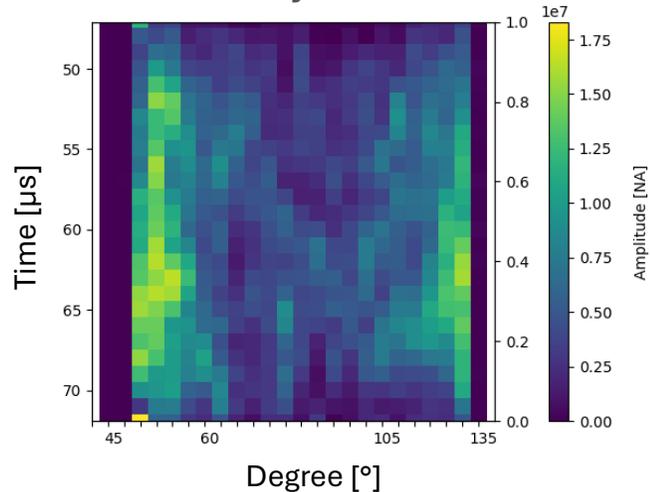# Rapid-E sampler

- Ranges:
  - 350-400 nm
  - 420-460 nm
  - 511-572 nm
  - 672-800 nm


Dactylus


Sensor Rapid-E
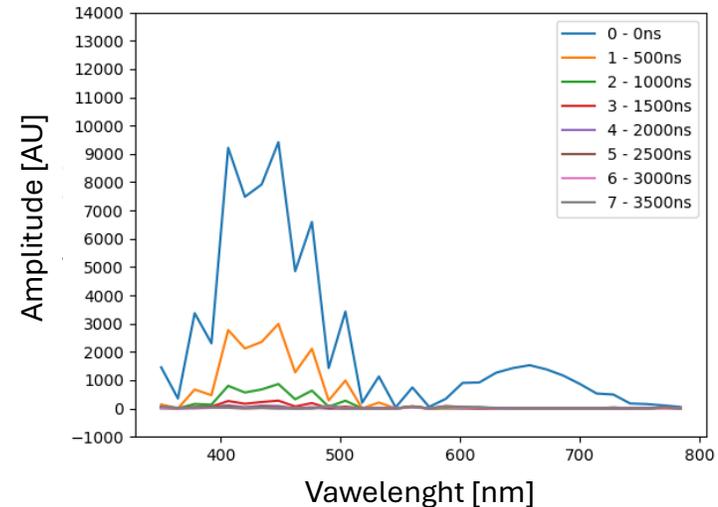
## Scattering


Corylus

## Fluorescence lifetime

Degree:
- 45-135 °

## Spectral ranges
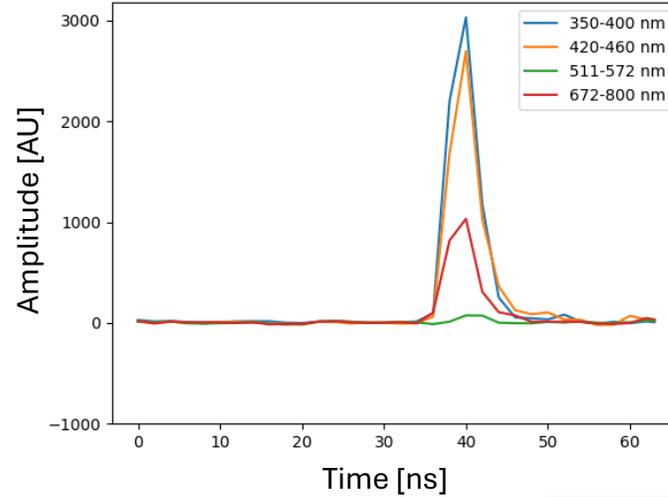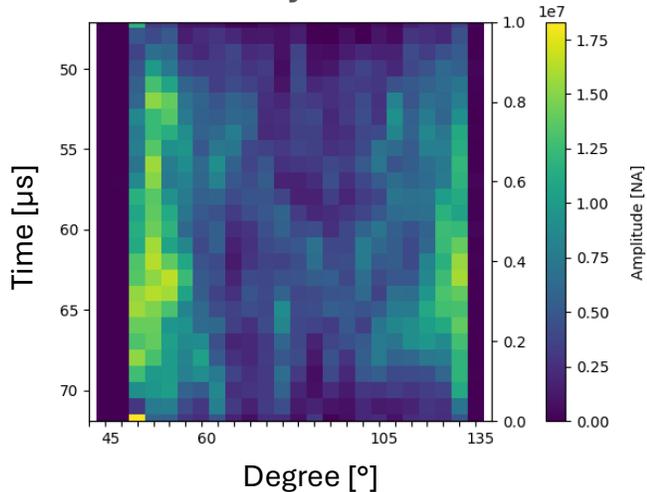

Alnus

# Rapid-E sampler

- Ranges:
  - 350-400 nm
  - 420-460 nm
  - 511-572 nm
  - 672-800 nm

**Dactylus**



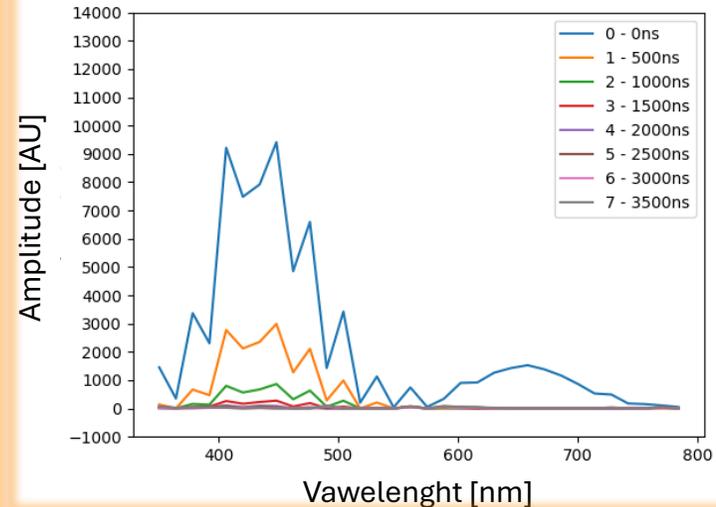**Sensor Rapid-E**

## Fluorescence lifetime

Degree:
- 45-135 °

Sampling
- 8 x 500 ns

## Scattering

**Corylus**



## Spectral ranges

**Alnus**

## Pollen type:

| | | | | | |
|---|---|---|---|---|---|
| 1. | Alnus | 6. | Fraxinus | 11. | Platanus |
| 2. | Arrhenatherum | 7. | Juglans | 12. | Populus alba |
| 3. | Broussonetia | 8. | Lolium perenne | 13. | Quercus |
| 4. | Corylus | 9. | Morus | 14. | Taxus |
| 5. | Dactylus | 10. | Pinus nigra | | |

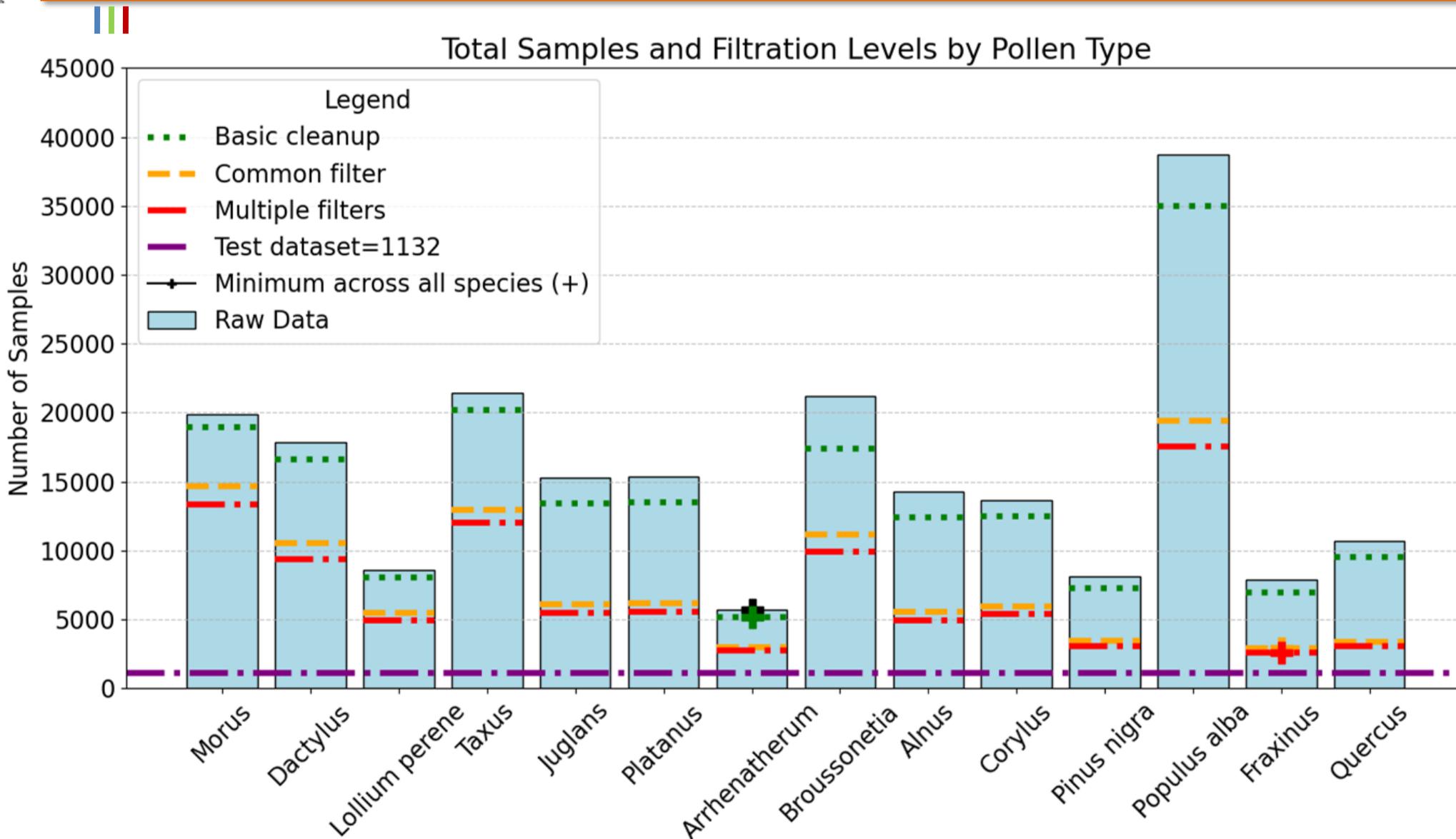*Plants surrounding INOE, collected by Boldeanu M.*

## 4 data filtration types selection

1. Raw data
2. Basic cleanup
3. One common filter
4. Pollen types divided into filtration groups

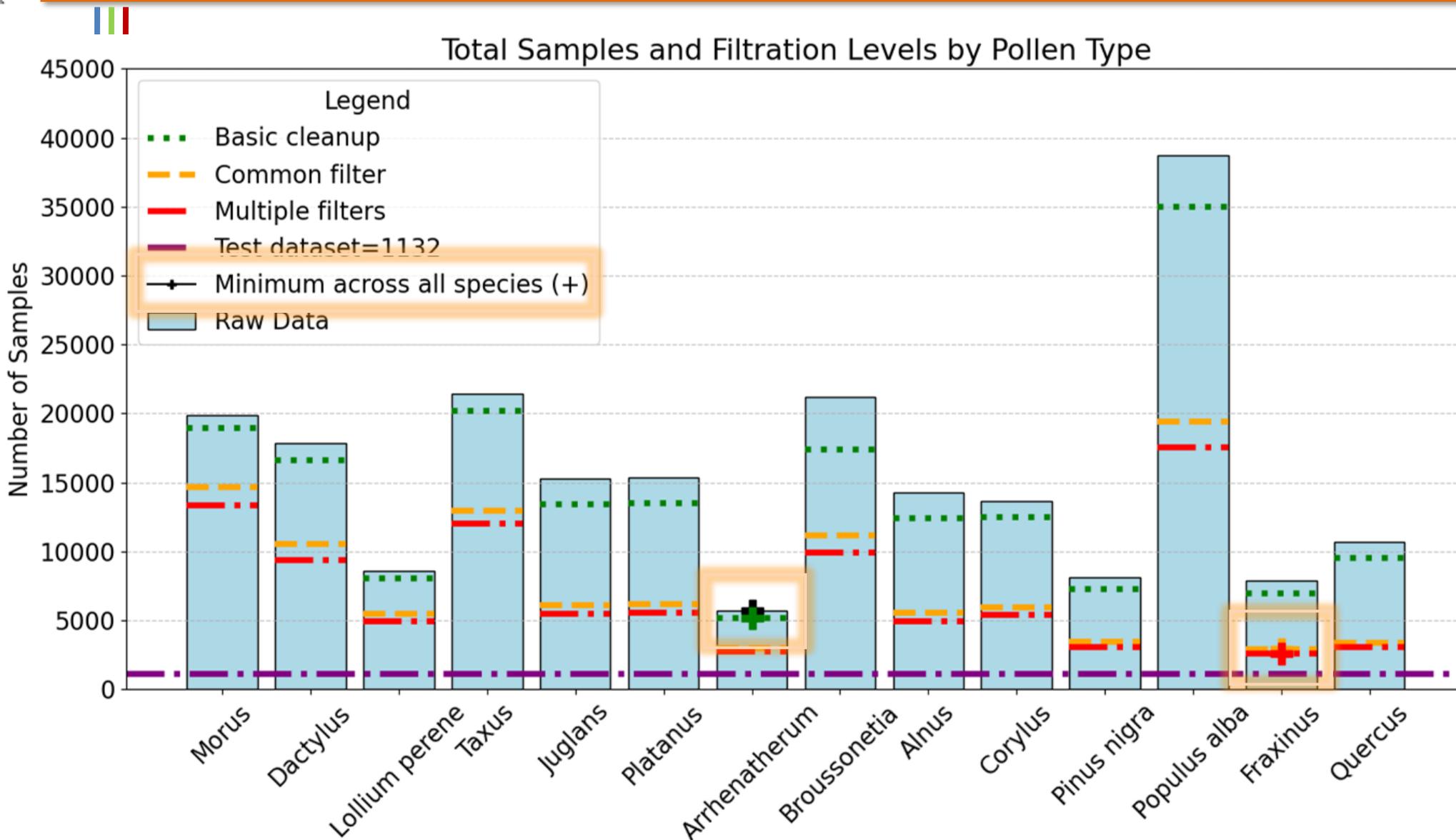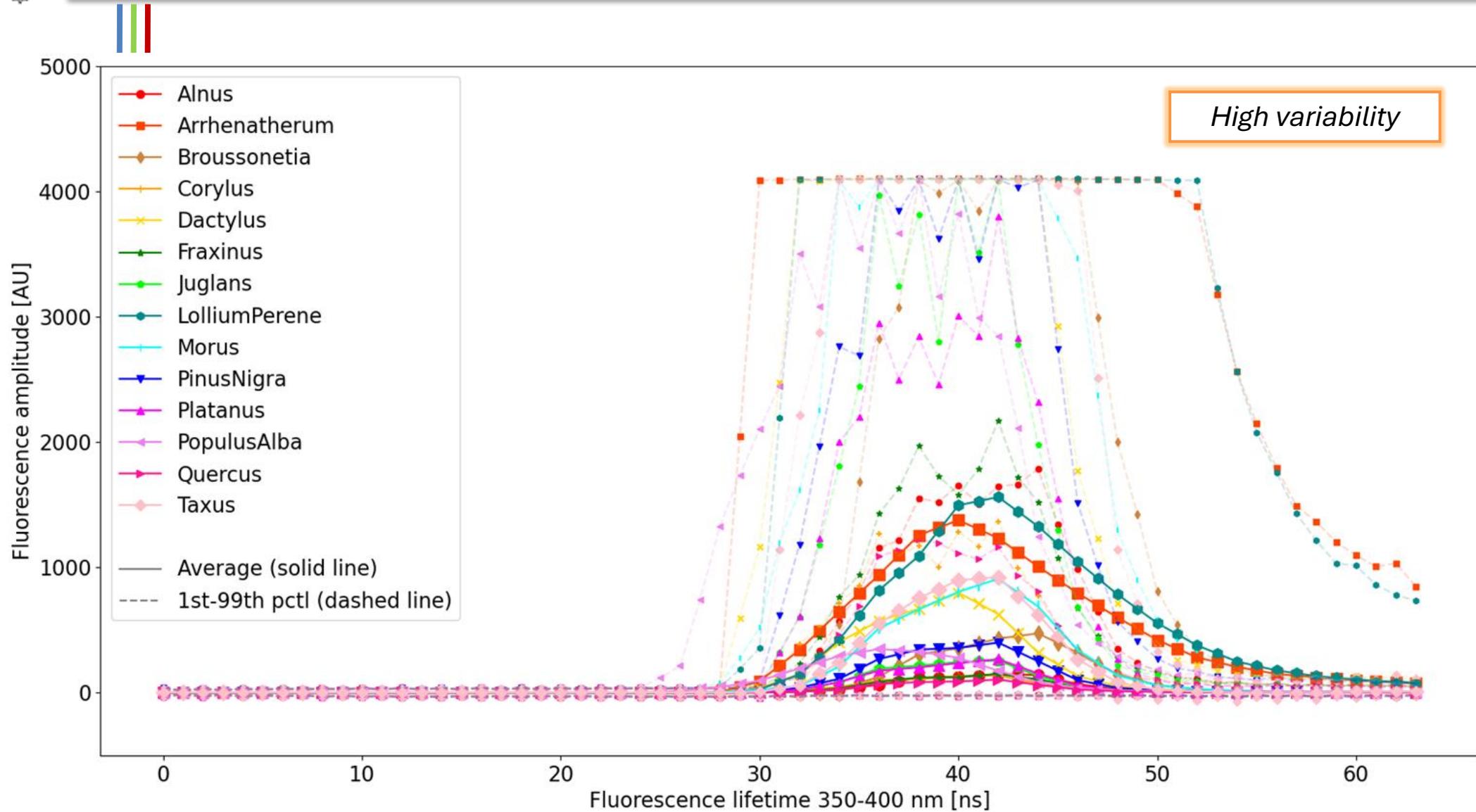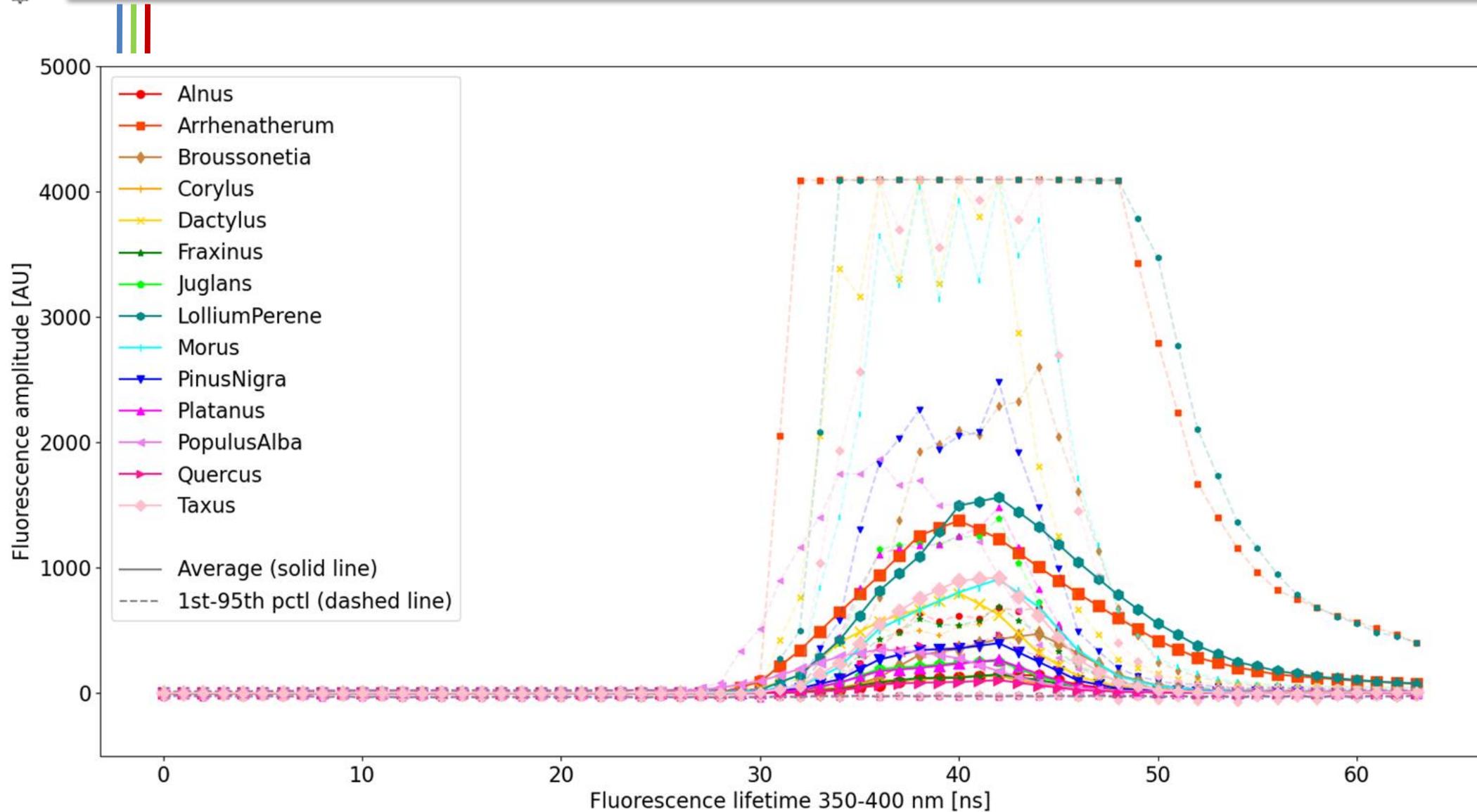*Same architecture for all models*

Total Samples and Filtration Levels by Pollen Type

Total Samples and Filtration Levels by Pollen Type

Spectrum after - 0 ns

Spectrum after - 0 ns

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS



*Model accuracy based on different filtration types. An average after 5 trials. Raw data represents the total number of samples*

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS



Model accuracy based on different filtration types. An average after 5 trials. Raw data represents the total number of samples
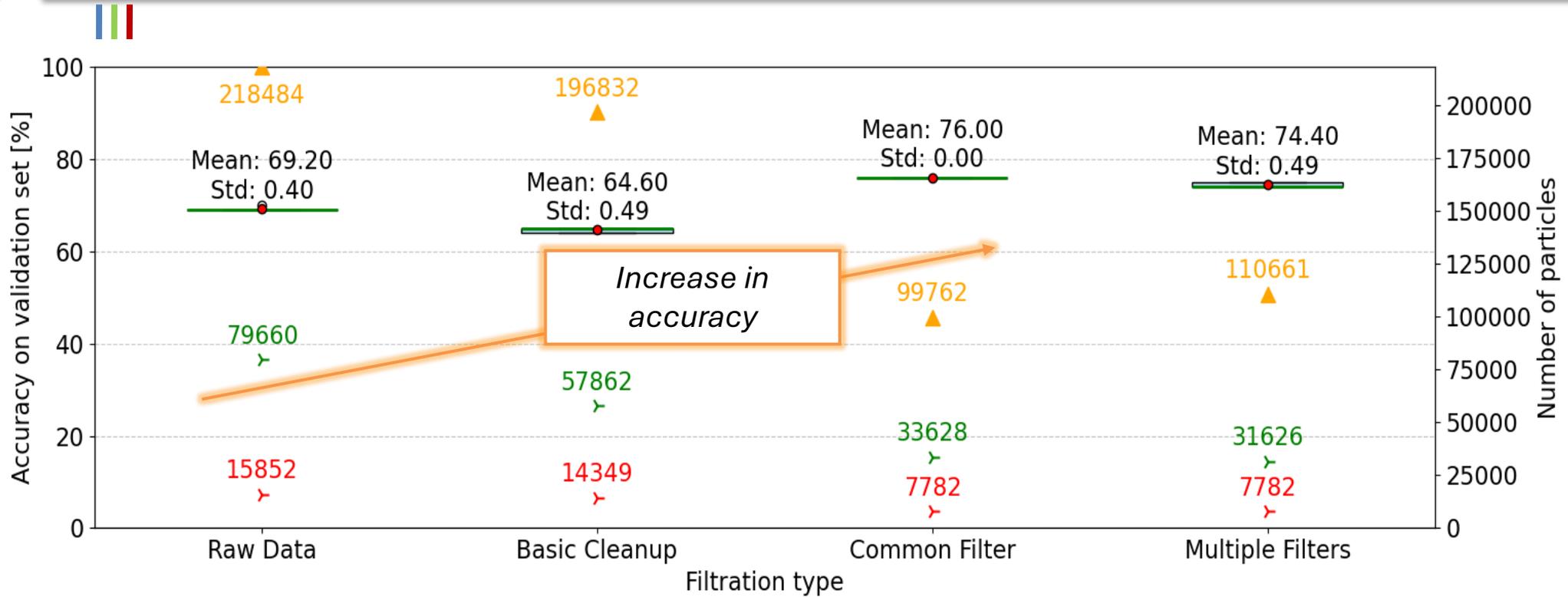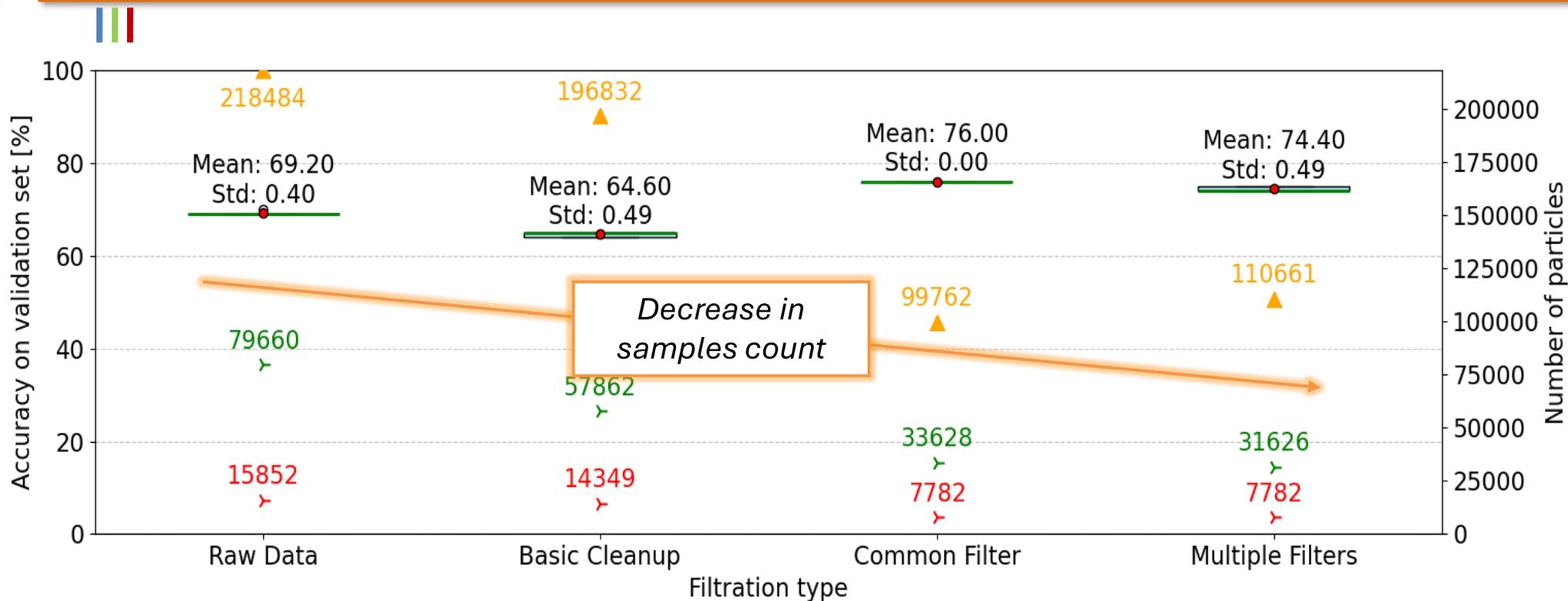
*Model accuracy based on different filtration types. An average after 5 trials. Raw data represents the total number of samples*

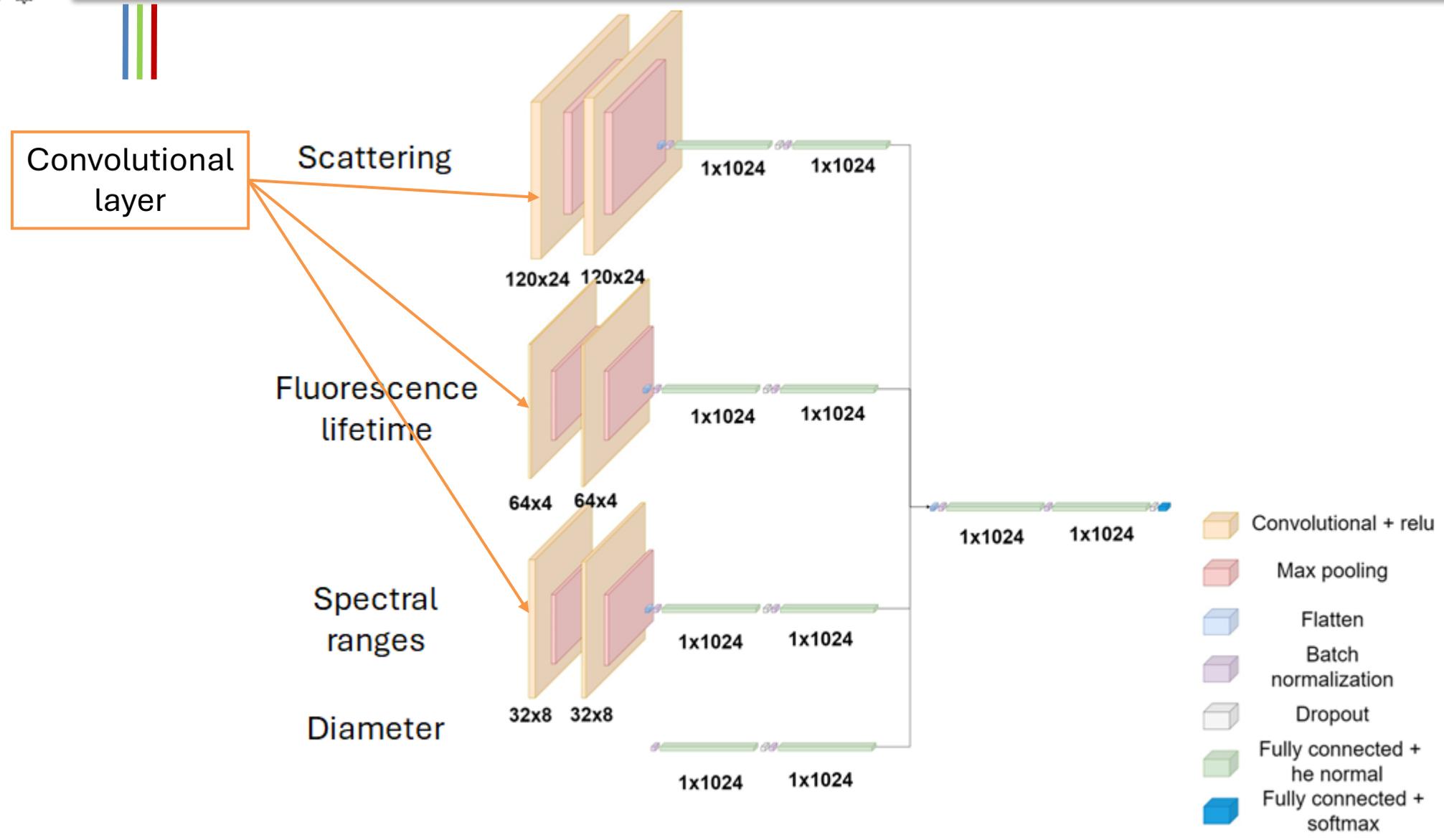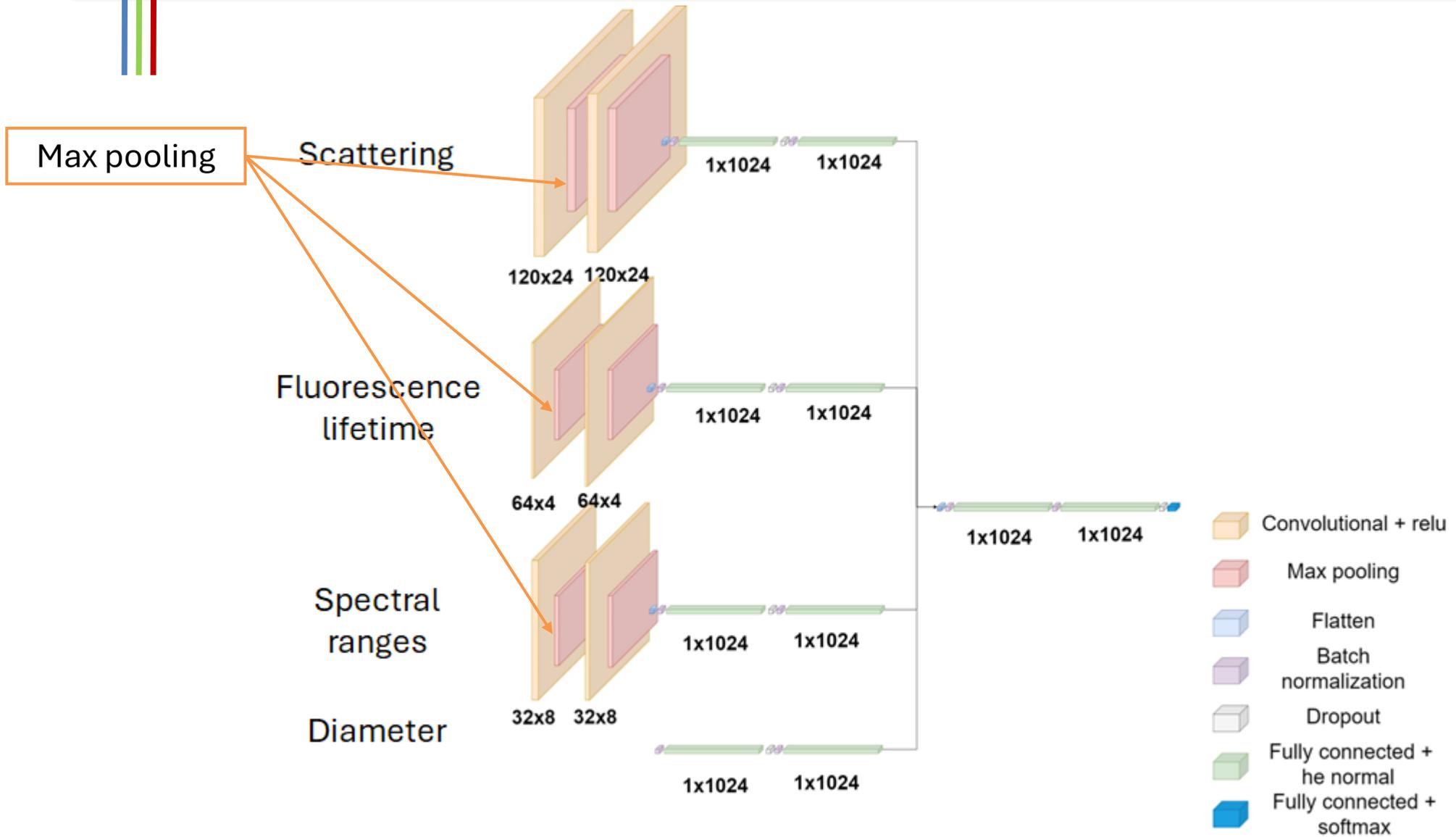UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS



Convolutional layer

Scattering
120x24   120x24
1x1024   1x1024

Fluorescence lifetime
64x4   64x4
1x1024   1x1024

Spectral ranges
32x8   32x8
1x1024   1x1024

Diameter
1x1024   1x1024

1x1024   1x1024
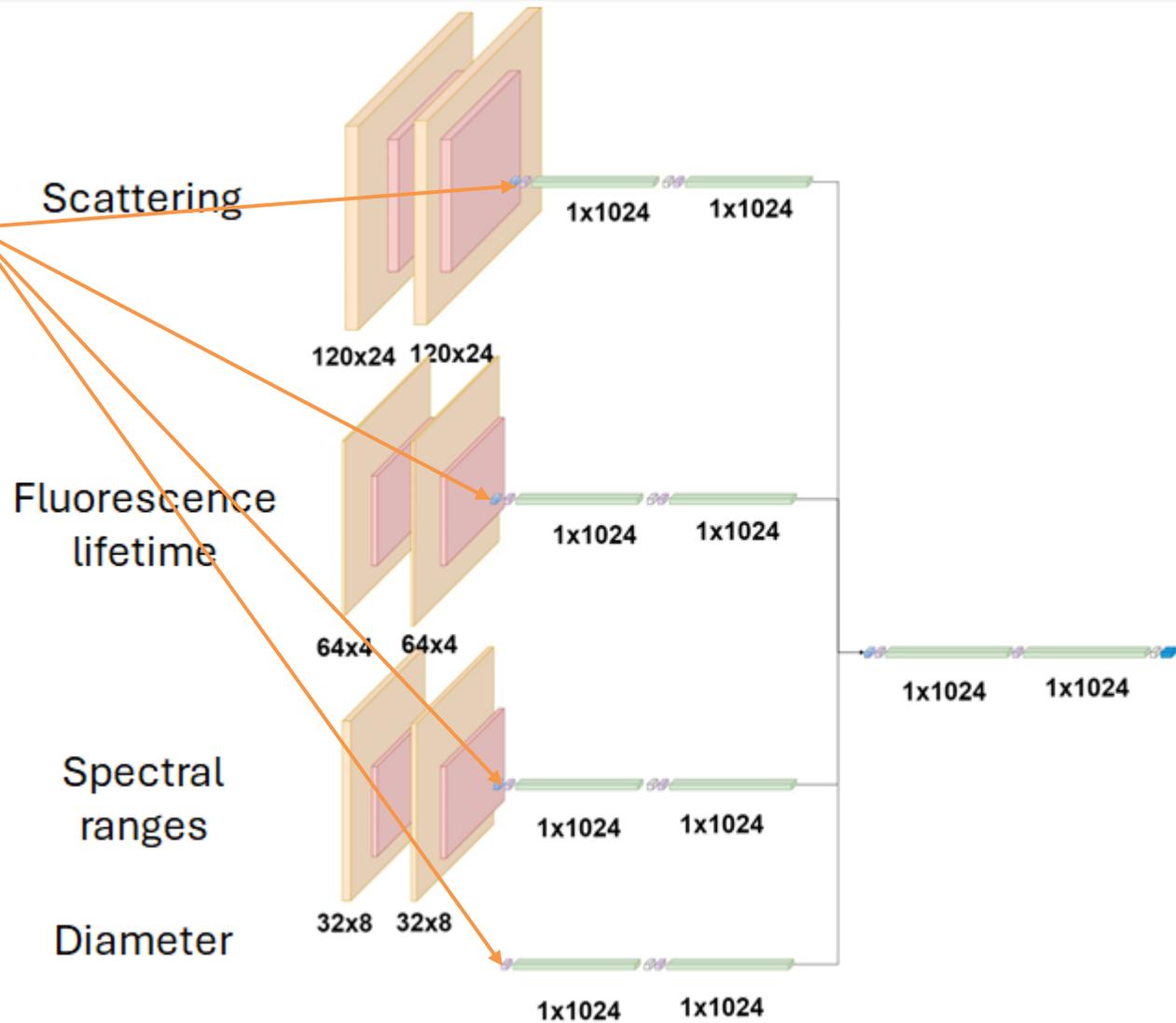
Convolutional + relu
Max pooling
Flatten
Batch normalization
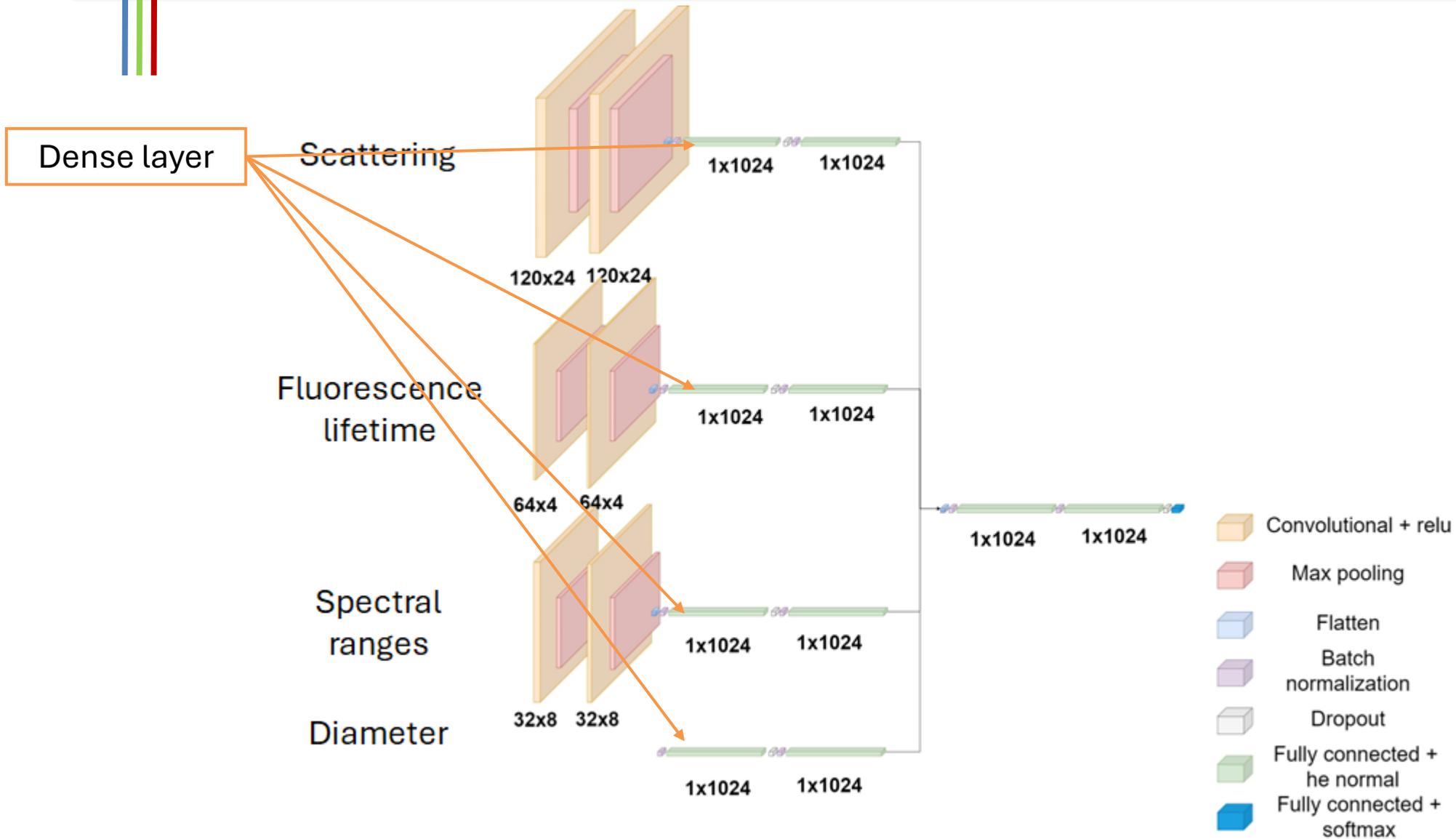Dropout
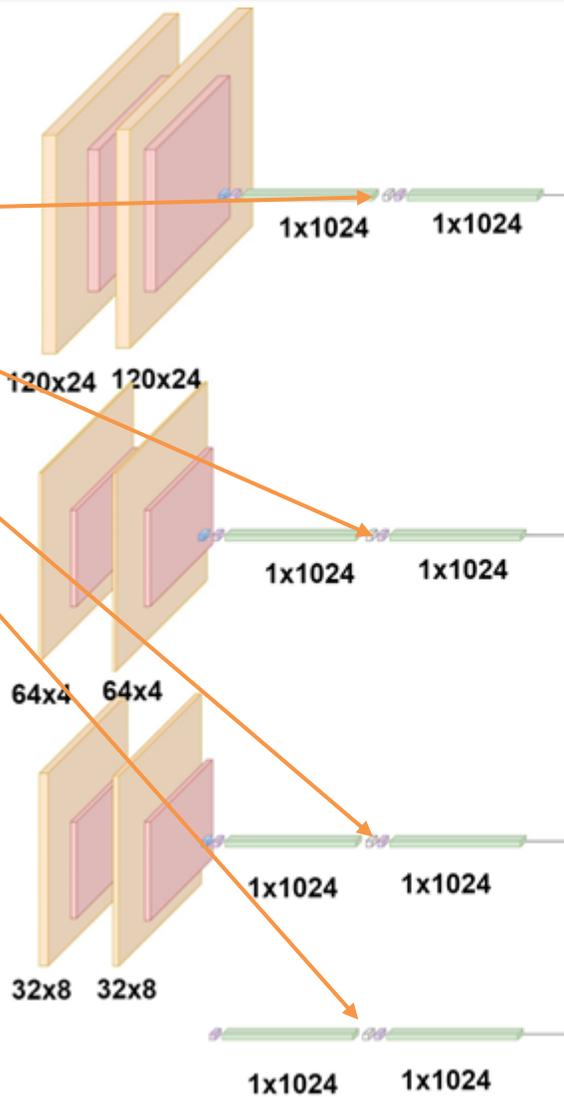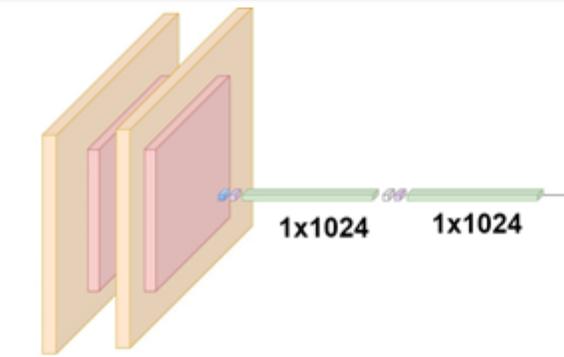Fully connected + he normal
Fully connected + softmax

Comparison of:

1. Precision and recall vs threshold
2. Total accuracy
3. Identification count vs threshold
4. ROC curve
5. Accuracy vs threshold
6. Confusion matrix

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

## Comparison of:

1. Precision and recall vs threshold
2. Total accuracy
3. Identification count vs threshold
4. ROC curve
5. Accuracy vs threshold
6. Confusion matrix

$T_p$ − true positived
$T_n$ − true negatives
$F_p$ − false positives
$F_n$ − false negatives
N − samples count
$p_o$ − threshold

$$f(p_o, precission, recall)$$

$$\text{precission} = \frac{T_p}{T_P + F_n}$$

$$recall = \frac{T_p}{T_p + F_p}$$

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

$T_l - true\ label$
$p_l - predicted\ label$

$$f(pl, Tp)$$

## Comparison of:

1. Precision and recall vs threshold
2. Total accuracy
3. Identification count vs threshold
4. ROC curve
5. Accuracy vs threshold
6. Confusion matrix

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

The higher precision and recall with lower threshold the better



Precision/recall vs threshold:
- recall
- precision
- compromise

Precision: 6/8 = 75%    4/5 = 80%    3/3 = 100%
Recall:    6/6 = 100%   4/6 = 67%    3/6 = 50%

Negative predictions — Various thresholds — Positive predictions — Score

~ Géron A., 2022

*Precision and recall – **raw data***

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS



*Precision and recall – **common filtered data***

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS



*Precision and recall – **raw data** (grey) vs **common filtered data** (coloured)*

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

The higher diagonal the better

The higher diagonal the better

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

The higher diagonal the better



| Common filter | 76% |
| Multiple filters | 75% |
| E-Rapid Clean | 65% |
| Raw data | 69% |

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

*Particle size histogram for training data – common filter*

Particle size histogram with Fraxinus and Taxus for training data – common filter

Overlapping sizes

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

*Particle size histogram for all pollen types for training data – common filter*

Overlapping sizes

# Case study

*One week of measurements comparison: 2022-05-10 – 2022-05-17*

*Multiple filters*

*Raw data*



*Common filter*

*Rapid-E cleanup*

*One week of measurements comparison: 2022-05-10 – 2022-05-17*

*Multiple filters*

*Raw data*

Daily fluctuations

*Common filter*

*Rapid-E cleanup*

*One week of measurements comparison: 2022-05-10 – 2022-05-17*

Common filter

Raw data

*Acceptable: 13/14*

*Acceptable: 8/14*

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

- Ready to use tool for:

  ~7k lines of code

  - Data filtering and inspecting - *DataViewer*
  - Training and validating models - *ModelBuilder*
  - Mapping results - *ModelRunner*
  - Visualising results - *PredictionsMapper*

- Ability to run in near real-time and map historical data

# Design assumptions

UNIVERSITY OF WARSAW | INSTITUTE OF GEOPHYSICS

- **User friendly**
  - User Interface (UI)
  - Config files
  - Logging

- **As small coding as possible on the end user**
  - Sometimes coding is necessary
  - Sometimes its easier to do something via code

- **Scalability**
  - Ability to easily add a new features
  - Large files processing

- **Adaptability**
  - Checkpoints
  - Ability to process data from different instruments

# Config files

# Logging



2025-03-13 22:42:54,617 | INFO | Logging in: '/home/tomcz/PollenTypingWsl/Rapid-E/resources/modelBuilder/logs/20250313.log'
2025-03-13 22:42:54,618 | INFO | Python version='3.11.4 (main, Jul  5 2023, 13:45:01) [GCC 11.2.0]'. (Implemented in 3.11.4)
2025-03-13 22:42:54,618 | INFO | Current Python interpreter path: '/home/tomcz/miniconda3/bin/python'
2025-03-13 22:42:54,618 | INFO | Configs cache on save is NOT invalidated
2025-03-13 22:43:00,927 | INFO | NumExpr defaulting to 16 threads.
2025-03-13 22:43:02,229 | INFO | Tensorflow version='2.12.1'. (Implemented in 2.12.1)
2025-03-13 22:43:02,732 | INFO | is_gpu_available: [PhysicalDevice(name='/physical_device:GPU:0', device_type='GPU')]
2025-03-13 22:43:02,732 | INFO | is_built_with_cuda=True
2025-03-13 22:43:03,965 | INFO | Cached pollen types file found. Retrieving. Filename='/home/tomcz/PollenTypingWsl/Rapid-E/resources/modelBuilder/cache/supervised/new_filter_thes
is/multiple_filters5000.h5'
2025-03-13 22:43:30,704 | INFO | Test file found. Retrieveing. File=/home/tomcz/PollenTypingWsl/Rapid-E/resources/modelBuilder/cache/supervised/new_filter_thesis/raw_data_test_mo
del.h5
2025-03-13 22:43:46,785 | INFO | Filtering test samples from the dataset, ones that does not fit restrictions.
2025-03-13 22:43:50,660 | INFO | Filtered samples count=8094/11922
2025-03-13 22:43:51,266 | INFO | Veryfying test set leaks skipped. To turn it on set 'veryfi_test_set_leaks' flag to 'true'
2025-03-13 22:43:54,562 | WARNING | Training not run due to proeprty 'run_training: false'. Change it to 'true' to process with the training.
2025-03-13 22:43:54,563 | INFO | Loading model from path: /home/tomcz/PollenTypingWsl/Rapid-E/resources/modelBuilder/multiple_filters5000_model.h5

```
2025-03-13 22:10:48,128 | INFO | Started - model runner
2025-03-13 22:10:48,128 | INFO | Particles to identify from files batch count=~100
2025-03-13 22:10:48,128 | INFO | Predictions threshold=0.8 (any pred for given particle above)
2025-03-13 22:10:48,128 | INFO | Output dir=modelRunner/out/thesis/multiple_filters_model
2025-03-13 22:10:49,913 | INFO | Running model: multiple_filters_model.h5
2025-03-13 22:10:49,923 | INFO | Progress count restored form path=/home/tomcz/PollenTypingWsl/Rapid-E/resources/modelRunner/out/thesis/multiple_filters_model/proc
2025-03-13 22:10:49,923 | INFO | Particles: identified_tr_0.8=64277, identified_tr_0=736426, total=2177513251
2025-03-13 22:10:49,924 | INFO | Searching data under path=/mnt/e/doktorat/raw_2022/d_00155
2025-03-13 22:10:49,990 | INFO | Searching data under path=/mnt/e/doktorat/raw_2022/d_00156
2025-03-13 22:10:50,057 | INFO | Searching data under path=/mnt/e/doktorat/raw_2022/d_00157
2025-03-13 22:10:50,227 | INFO | Searching data under path=/mnt/e/doktorat/raw_2022/d_00158
2025-03-13 22:10:50,295 | INFO | Searching data under path=/mnt/e/doktorat/raw_2022/d_00159
2025-03-13 22:10:51,182 | INFO | Identified dates: dict_keys(['20220426', '20220427', '20220428', '20220429', '20220430', '20220501', '20220502', '20220503', '2022
2025-03-13 22:10:51,212 | INFO | All files for date=20220426 were processed. Date will be skipped.
2025-03-13 22:10:51,260 | INFO | All files for date=20220427 were processed. Date will be skipped.
2025-03-13 22:10:51,328 | INFO | All files for date=20220428 were processed. Date will be skipped.
2025-03-13 22:10:51,421 | INFO | All files for date=20220429 were processed. Date will be skipped.
2025-03-13 22:10:59,742 | INFO | Files to process count=24970
2025-03-13 22:10:59,743 | INFO | Processing date = 2022-04-26 00:00:00 (1/35)
2025-03-13 22:10:59,744 | INFO | Processing date = 2022-04-27 00:00:00 (2/35)
2025-03-13 22:10:59,744 | INFO | Processing date = 2022-04-28 00:00:00 (3/35)
2025-03-13 22:10:59,744 | INFO | Processing date = 2022-04-29 00:00:00 (4/35)
2025-03-13 22:11:01,898 | INFO | Particles: identified_tr_0.8=64278, identified_tr_0=736531, total=2177513356
2025-03-13 22:11:01,902 | INFO | Processed files count=3/24970
2025-03-13 22:11:02,311 | INFO | Particles: identified_tr_0.8=64280, identified_tr_0=736647, total=2177513577
2025-03-13 22:11:02,314 | INFO | Processed files count=8/24970
2025-03-13 22:11:02,677 | INFO | Particles: identified_tr_0.8=64280, identified_tr_0=736754, total=2177513905
2025-03-13 22:11:02,680 | INFO | Processed files count=11/24970
```
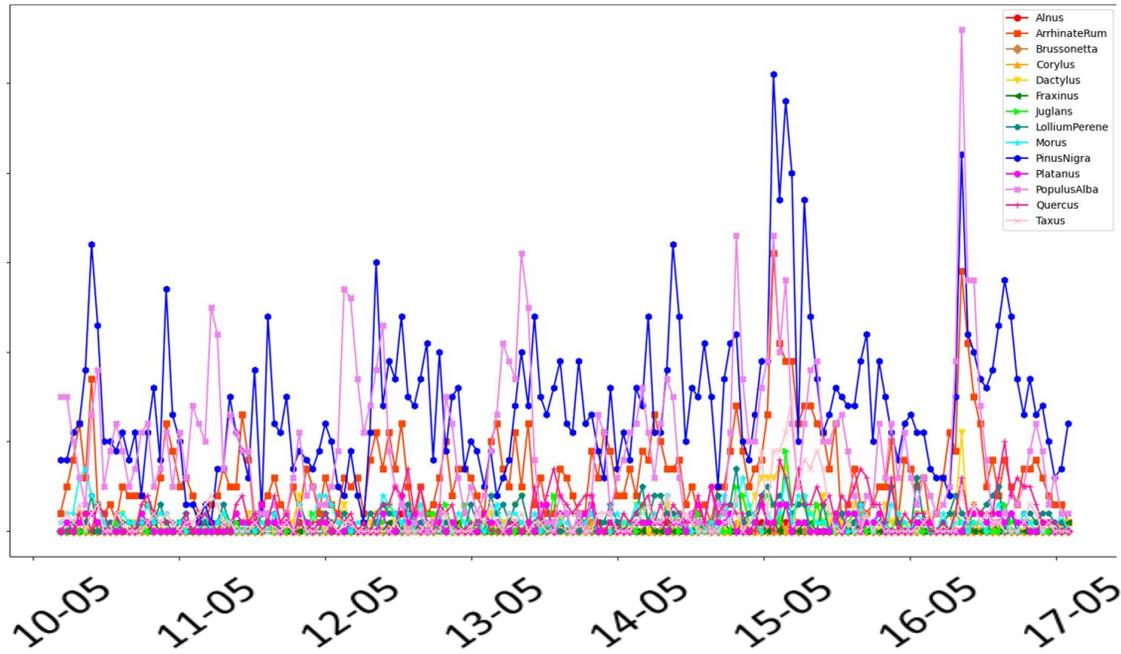
- ## Current status:
  - ModelRunner and PredictionsMapper – fully independent
  - DataViewer and ModelBuilder – partially independent – requires some adaptations
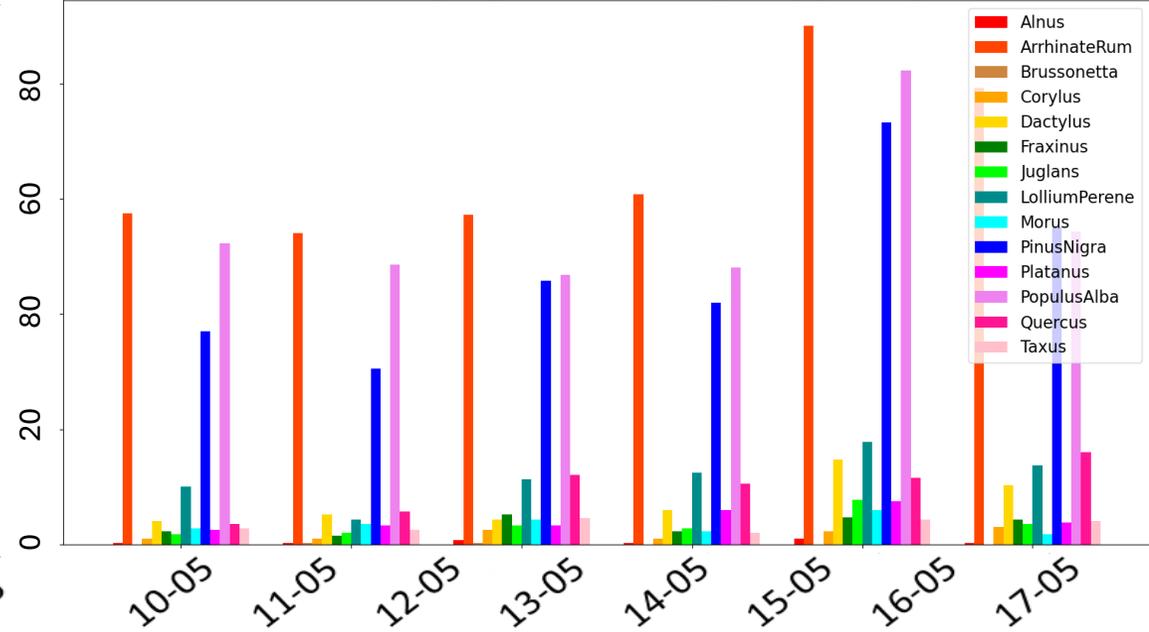
- ## Ability to read:
  - Plair binary files – raw data files from instrument
  - .json files

# Visualising the data

- Promising results – validation required

- Ready to use tool for:
  - Data filtering and inspecting
  - Training models
  - Validating models
  - Mapping results
  - Visualising results

- Ability to run in near real-time and map historical data

# Thank you

Contact emails:        artur.tomczak@fuw.edu.pl
                       iwona.stachlewska@fuw.edu.pl