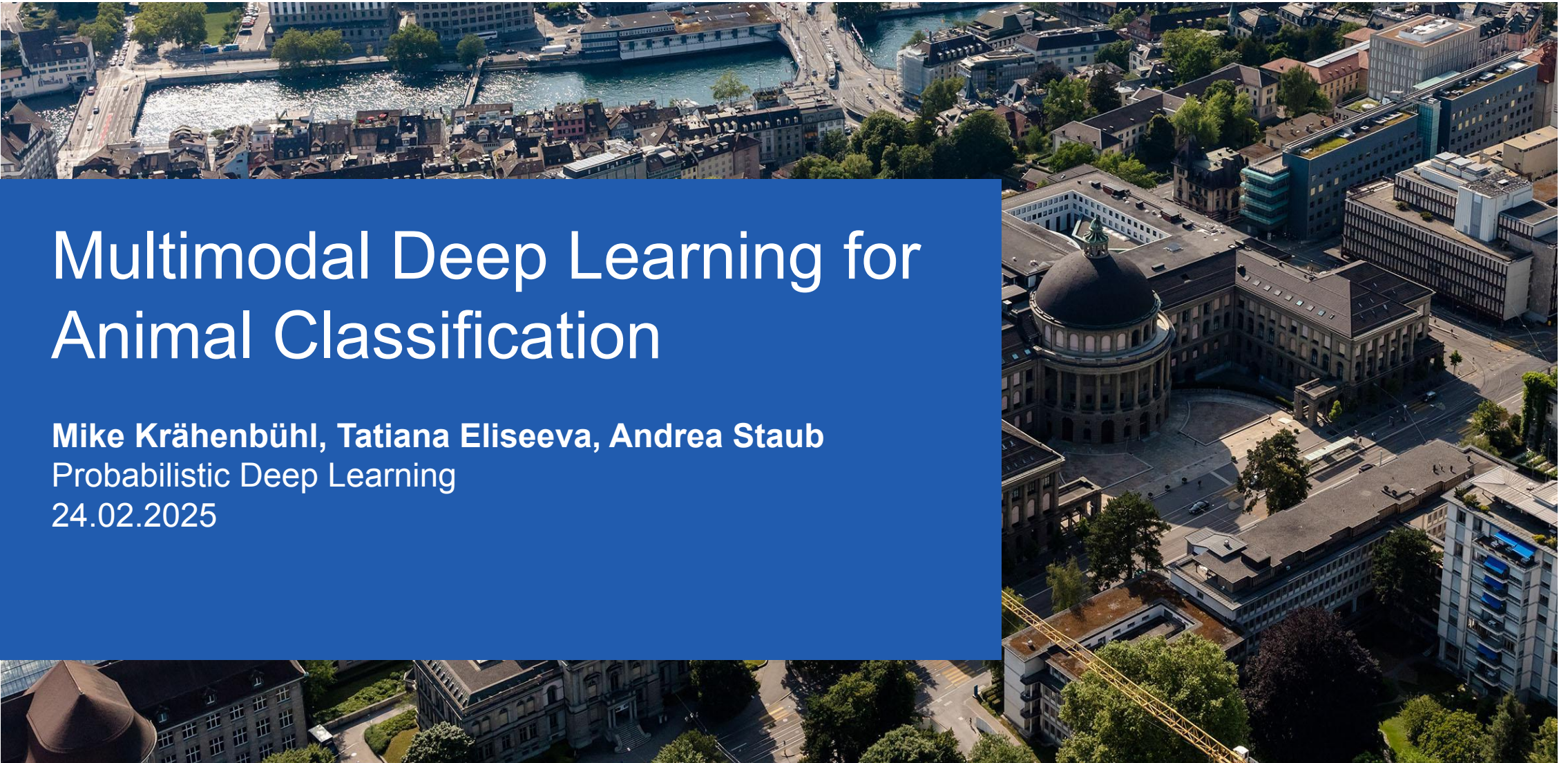


# Multimodal Deep Learning for Animal Classification

**Mike Krähenbühl, Tatiana Eliseeva, Andrea Staub**

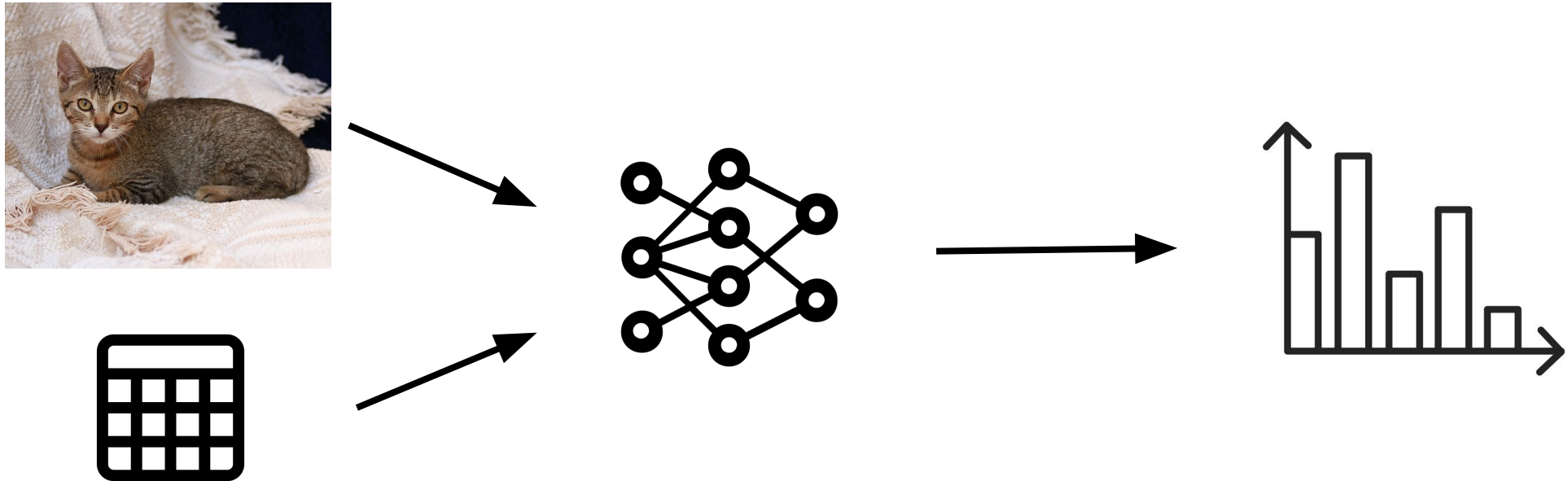
Probabilistic Deep Learning

24.02.2025



# Goal

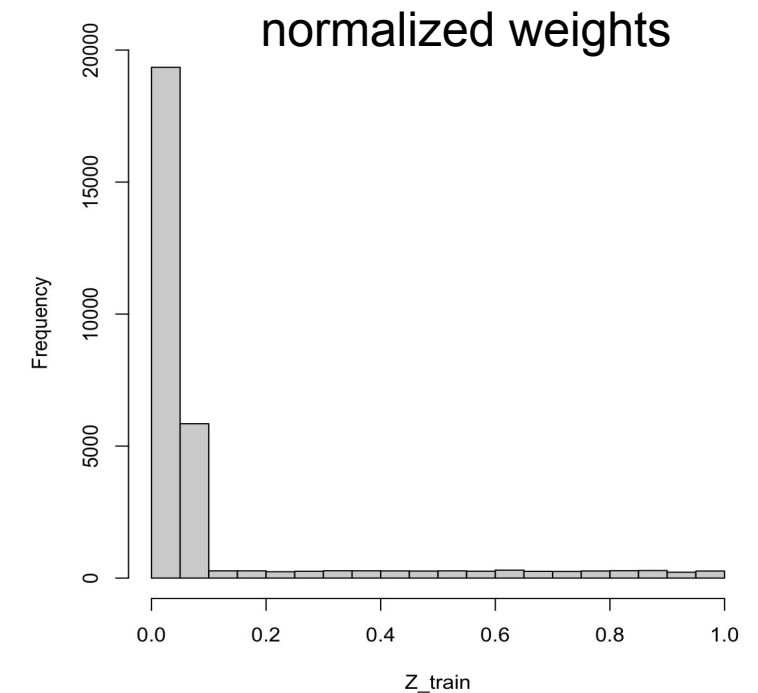
- Predicting the probability distribution of animal categories based on **image** data and **tabular** data using different NN architectures.





# Data

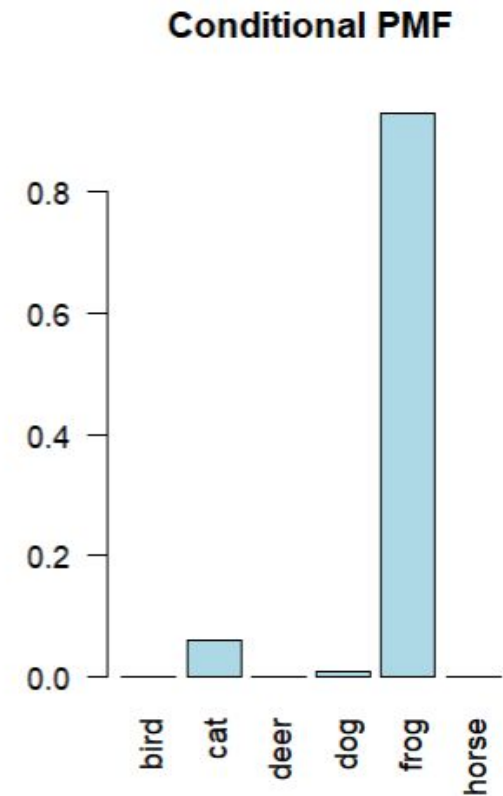
- **Image data:** CIFAR-10 animal images, 32x32 colour images in six classes, with 6'000 images per class.
- **Tabular data:** Randomly generated weight from uniform distribution with reasonable lower and upper bounds for each animal type.



# Estimated Conditional Probability Distribution for an Observation

Model:

CNN (image) + Tabular (weight)

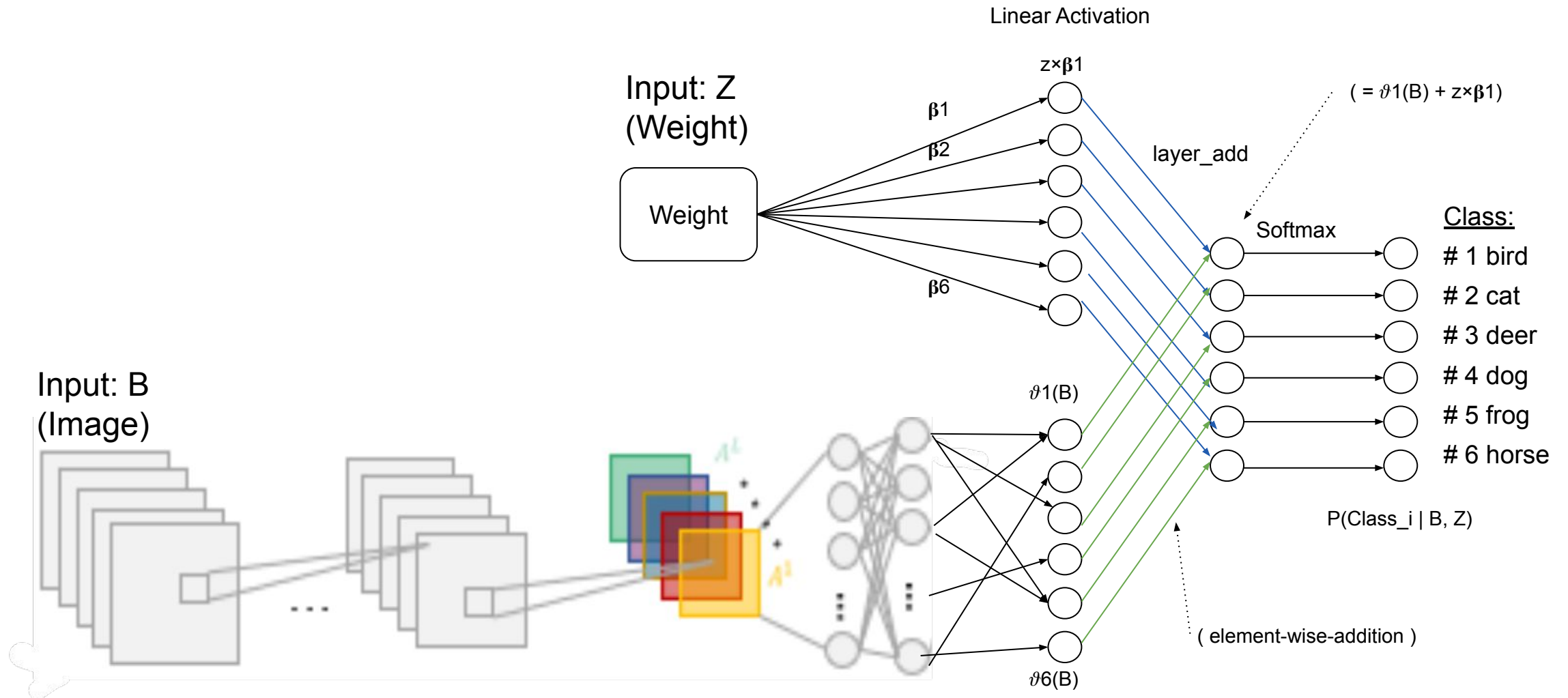


# Models to compare

| Model Name                | Architecture Brief Summary*  | Inputs             |
|---------------------------|--|--------------------|
| Model 1 - CNN             | “Classic” Image CNN  | Images             |
| Model 2 - 1LfcNN-b        | Single-Layer fcNN. Contains weights only, softmax activation function without bias | Weights            |
| Model 3 - 1LfcNN          | Single-Layer fcNN. Contains weights only, softmax activation function with bias    | Weights            |
| Model 4 - Deep-fcNN       | Deep fcNN. Contains weights only, ReLu kern with bias                              | Weights            |
| Model 5 - CNN + 1LfcNN    | CNN + Single-Layer fcNN  | Images and weights |
| Model 6 - CNN + Deep-fcNN | CNN + Deep fcNN  | Images and weights |

\*Always softmax activation at output nodes

# Model 5 Architecture (CNN + linearNN)



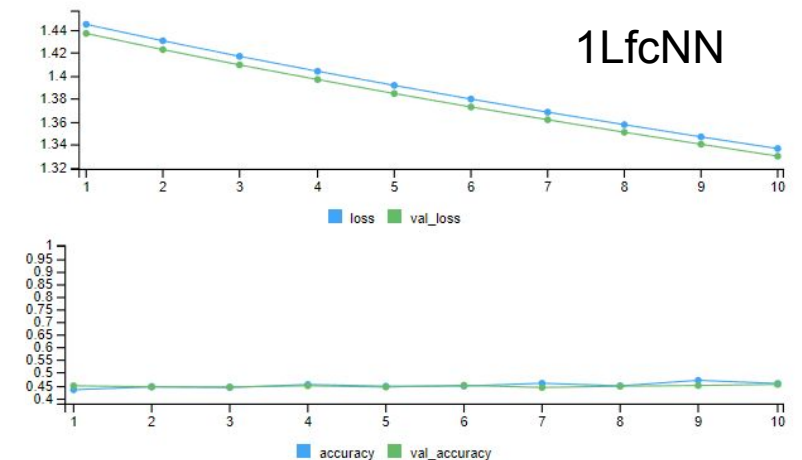
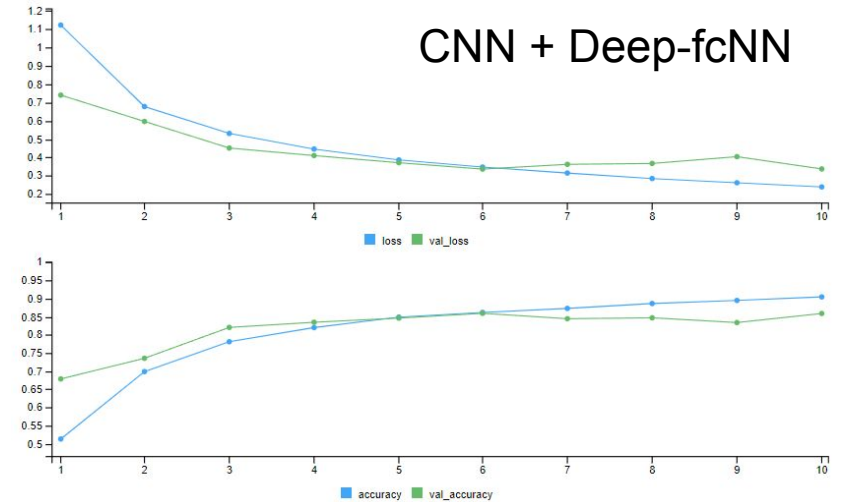
# Training Workflow

1. Input Normalisation and One-hot encoding for classes
2. Model compilation with NLL (loss = "loss\_categorical\_crossentropy", optimizer = optimizer\_adam(), metrics = c("accuracy"))
3. Initialize Early Stopping parameters (callback\_early\_stopping(monitor = "val\_loss", patience = 3, restore\_best\_weights = TRUE))
4. Model training: epochs = 100, batch\_size = 64, validation\_split = 0.2, callbacks = list(early\_stopping)

Multiple model initialisation did not make any significant difference for the model accuracy. As an example - validation loss about two digits after comma. It shows that the model is robust for the chosen data.

Deviations on “classical settings”, sometimes quite big, like using only half of the nodes, did not brought any significant results.

“Thumb Rule” for our models and our problems: 10 Epochs for the run is enough.

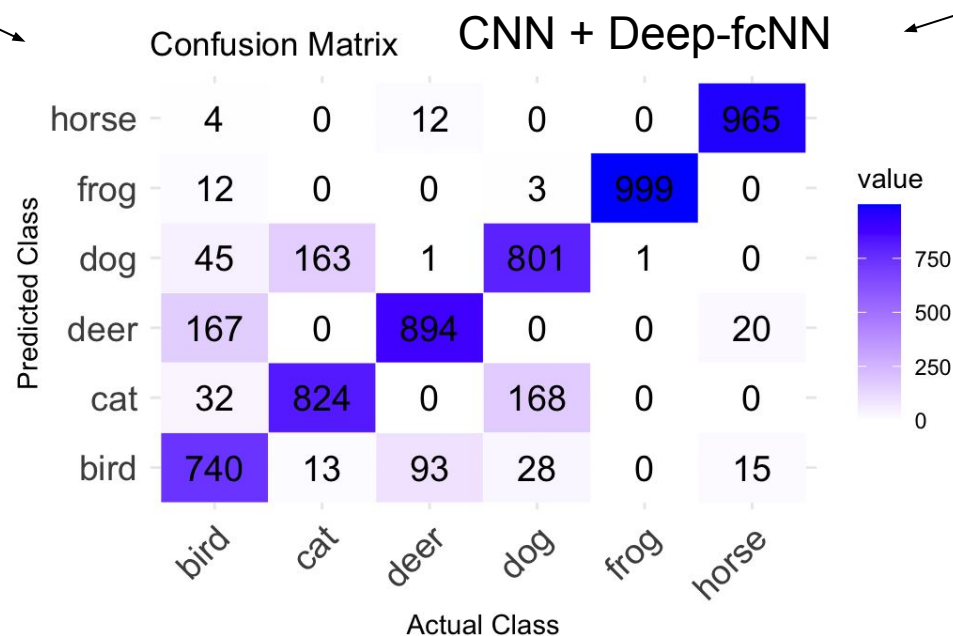
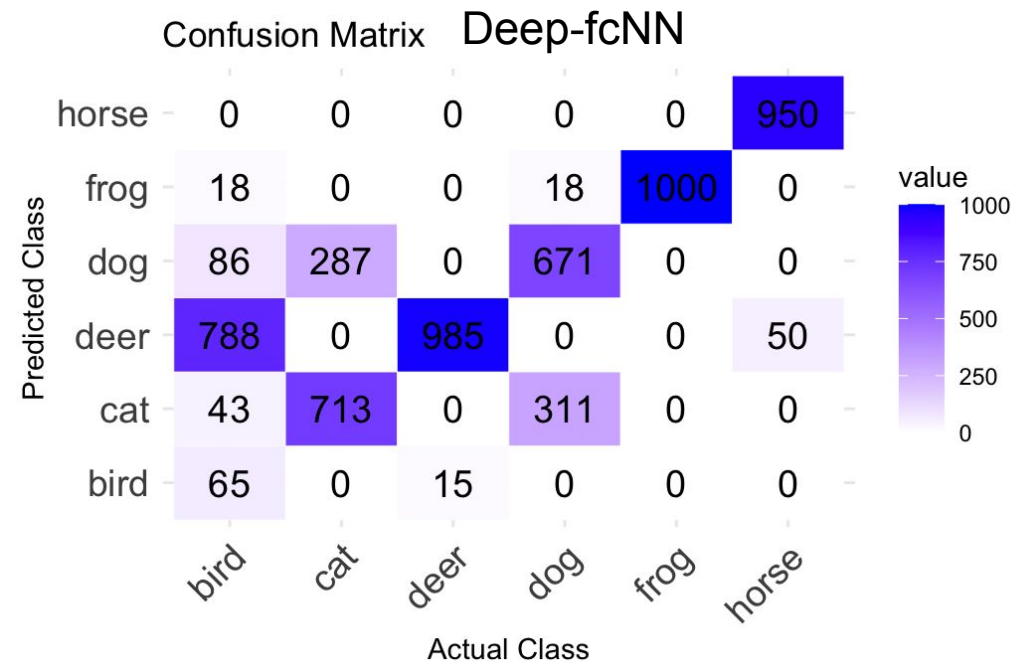
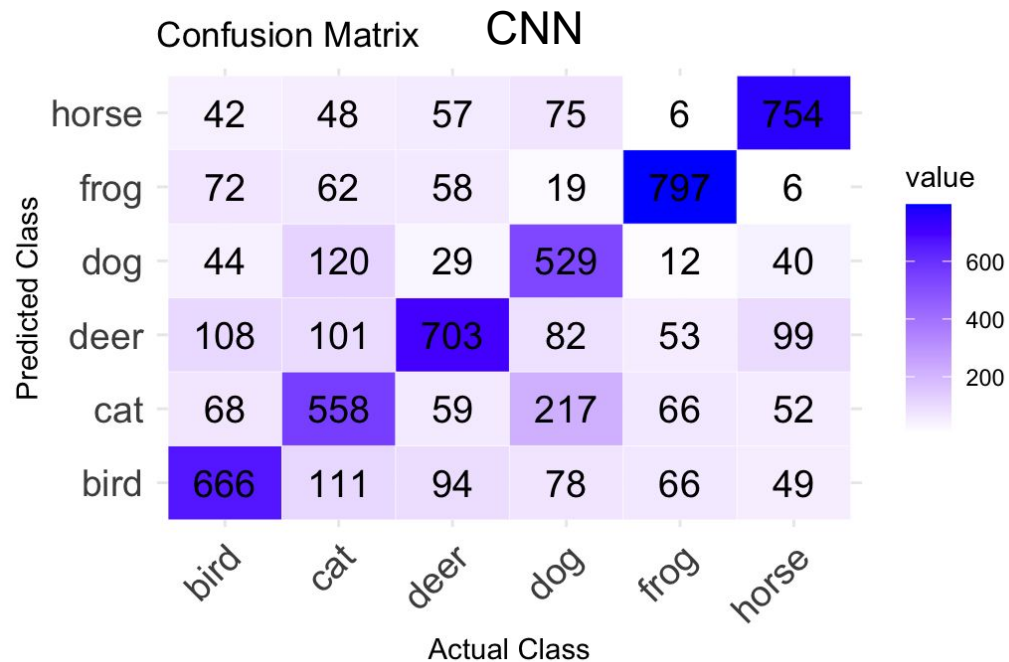


# Test Results

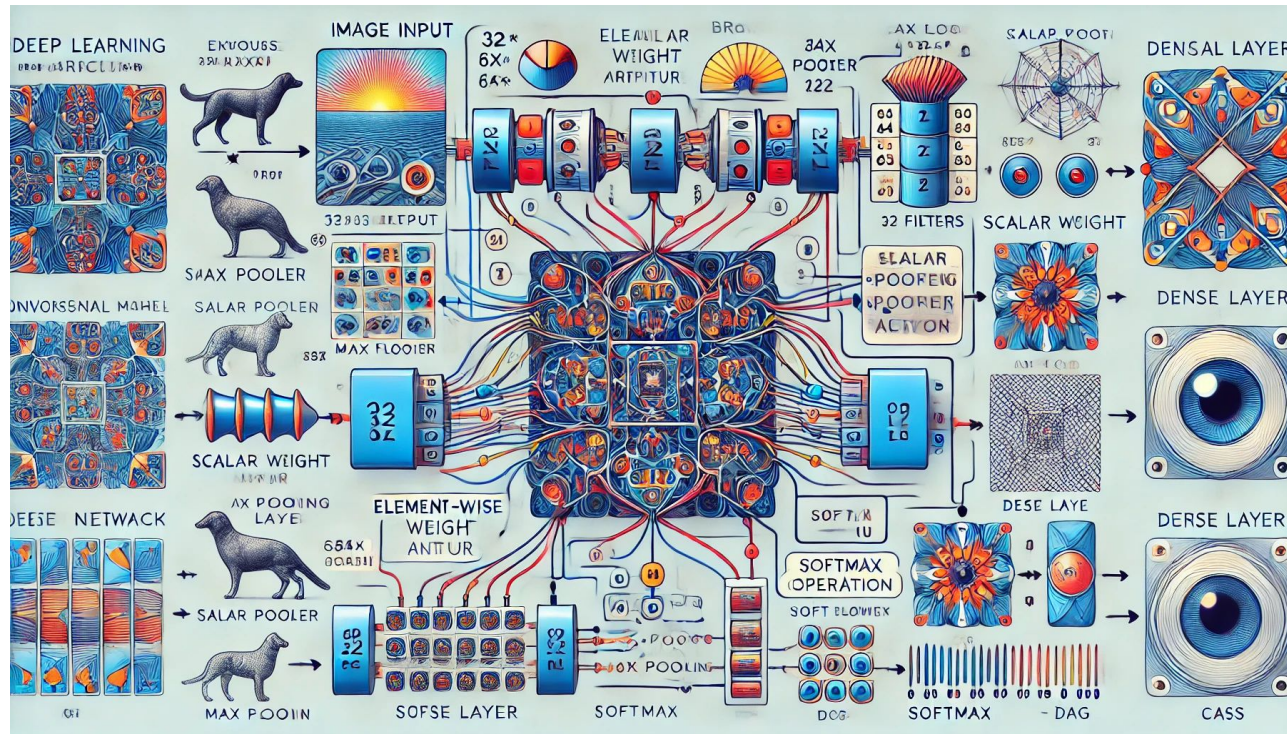
| Model           | Accuracy | Loss  |
|-----------------|----------|-------|
| CNN             | 0.674    | 0.902 |
| 1LfcNN-b        | 0.167    | 1.766 |
| 1LfcNN          | 0.446    | 1.451 |
| Deep-fcNN       | 0.757    | 0.537 |
| CNN + 1LfcNN    | 0.713    | 0.792 |
| CNN + Deep-fcNN | 0.871    | 0.344 |

- Single-Layer NN has too low accuracy
- Deep NN allow to capture the features needed for proper categorisation
- Already one layers of fcNN that allows to process tabular data, can significantly increase the accuracy
- Best results shows combination of CNN and Deep fcNN





# Conclusions



- Combination of different types of input data by training of DeepNN tools allows to increase the recognition efficiency for the low-weight solutions
- For data without significant preprocessing a NN has to have sufficient capacity to store the features