

# Data shortage? Let Diffusion fill the gaps!

## ImbalanceSD

### Project overview

This project aims to investigate possible improvements that synthetic data generation can bring to the unbalanced classification problem. By generating new data samples using diffusion models and integrating them in an intelligent way, I hope to improve the classifier models.

### Image generation

4 diffusion models were used to generate data:

- FLUX.1-dev
- FLUX.1-schnell
- StableDiffusion 3.5 Large
- StableDiffusion 3.5 Large Turbo

All generated images have good quality, but FLUX.1-dev has produced the best looking ones.

### Used prompt:

*f"{camera\_angle} of a {breed} cat {preposition} the {furniture}, {cat\_gaze}. The cat has realistic fur textures, intricate details, and sharp features, with soft lighting and a clear focus. The image has a shallow depth of field, emphasizing the cat in fine detail. 8k, cinematic, photorealistic*

### Performed experiments

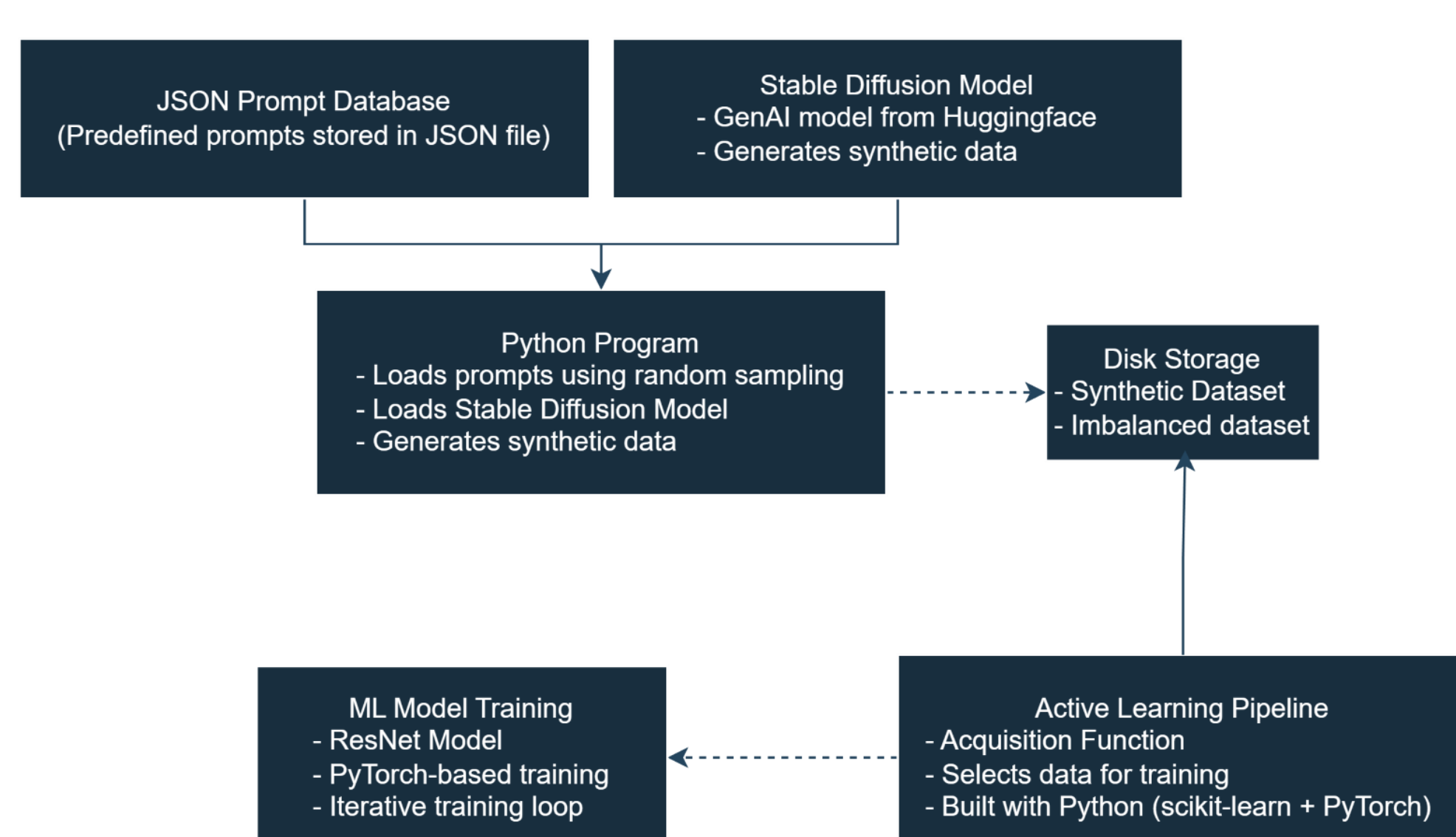
*Experiment (testset accuracy)*

- ResNet18 training on full CIFAR10 (86%)
- ResNet18 training on CIFAR10 with „cat“ class downsampled to 1% using:
  - Straightforward training (80%)
  - SMOTE, ADASYN and naive oversampling (81%)
  - Undersampling (21%)
  - Class weighting (76%)
  - Label smoothing (80%)
- Synthetic data from FLUX.1-dev (still in experiments, no results yet)

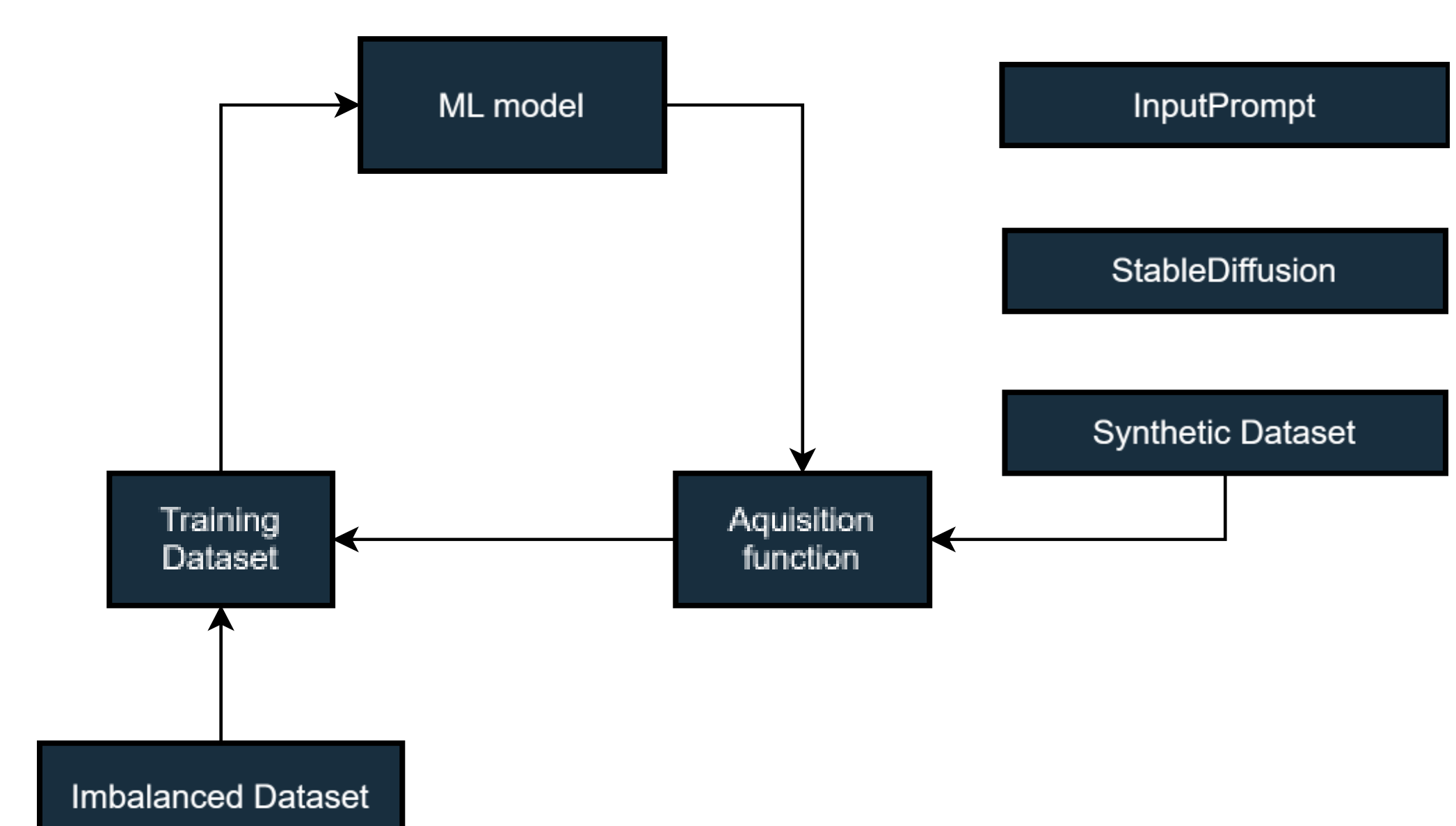
### Future steps

- Comparison of trainings using different synthetic data
- Usage of Active Learning's different acquisition function based on uncertainties
- Image augmentation using FLUX.1-Redux (diffusion model generating the same object from the photo, but slightly changed)
- More augmentations
- New images generation using FLUX.1-dev photorealism LoRA
- Model guided active generation during training

### System architecture diagram



### Active learning pipeline



Karol Cyganik  
Under the supervision of Dr Wojciech Pałubicki

