



GaNDLF-Synth: a Framework for Generative AI in Biomedical Imaging

Sarthak Pati^{1,2}, Szymon Mazurek^{3,4}, Akis Lindaros¹
Spyridon Bakas^{1,2}

1 Indiana University, Indianapolis, 46202 USA

2 Medical AI Working Group, MLCommons, San Francisco, CA, USA

3 AGH University of Krakow, Krakow, 30-059 Poland

4 ACC Cyfronet AGH, Krakow, 30-950 Poland

04.04.2025

Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie
AGH University of Krakow



Problem overview

- **Data Scarcity:** medical imaging data is limited compared to other data types
- **Privacy Concerns:** sharing patient medical images raises significant privacy issues
- **Data Imbalance:** rare diseases have very few positive cases, creating data imbalances
- **Inter-Site Variability:** medical images vary across different acquisition sites
- **Lack of General-Purpose Tools:** existing tools are often specialized and not easily adaptable
- **Computational Expertise Required:** training GenAI models demands significant technical knowledge
- **Translation Difficulties:** generating one modality of medical image from another is complex

GenAI and problems in medical imaging

Aspect	Description
Data Imbalance	GenAI can address gross data imbalance, especially for rare diseases, by generating more positive cases. It can also handle class imbalance across sites in multi-institutional studies Sheller et al. (2020) ; Pati et al. (2022) .
Unsupervised Learning	GenAI models can learn data distribution in an unsupervised manner and be transferred as feature extractors for other tasks, such as classification Lin et al. (2017) ; Shrivastava et al. (2017) .
Data Quality Enhancement	GenAI models can enhance data quality by performing tasks like denoising, reconstruction, or super-resolution Frangi et al. (2018) .
Modality Generation	With proper training, GenAI models can generate different imaging modalities from base examples, such as generating PET images from MR Dayarathna et al. (2023) , thus, serving as an approximation tool.
Privacy	GenAI can enable the use of synthesized data instead of real patient data, which can increase privacy by reducing the risk of models leaking training information Song et al. (2019) ; Pati et al. (2024) .



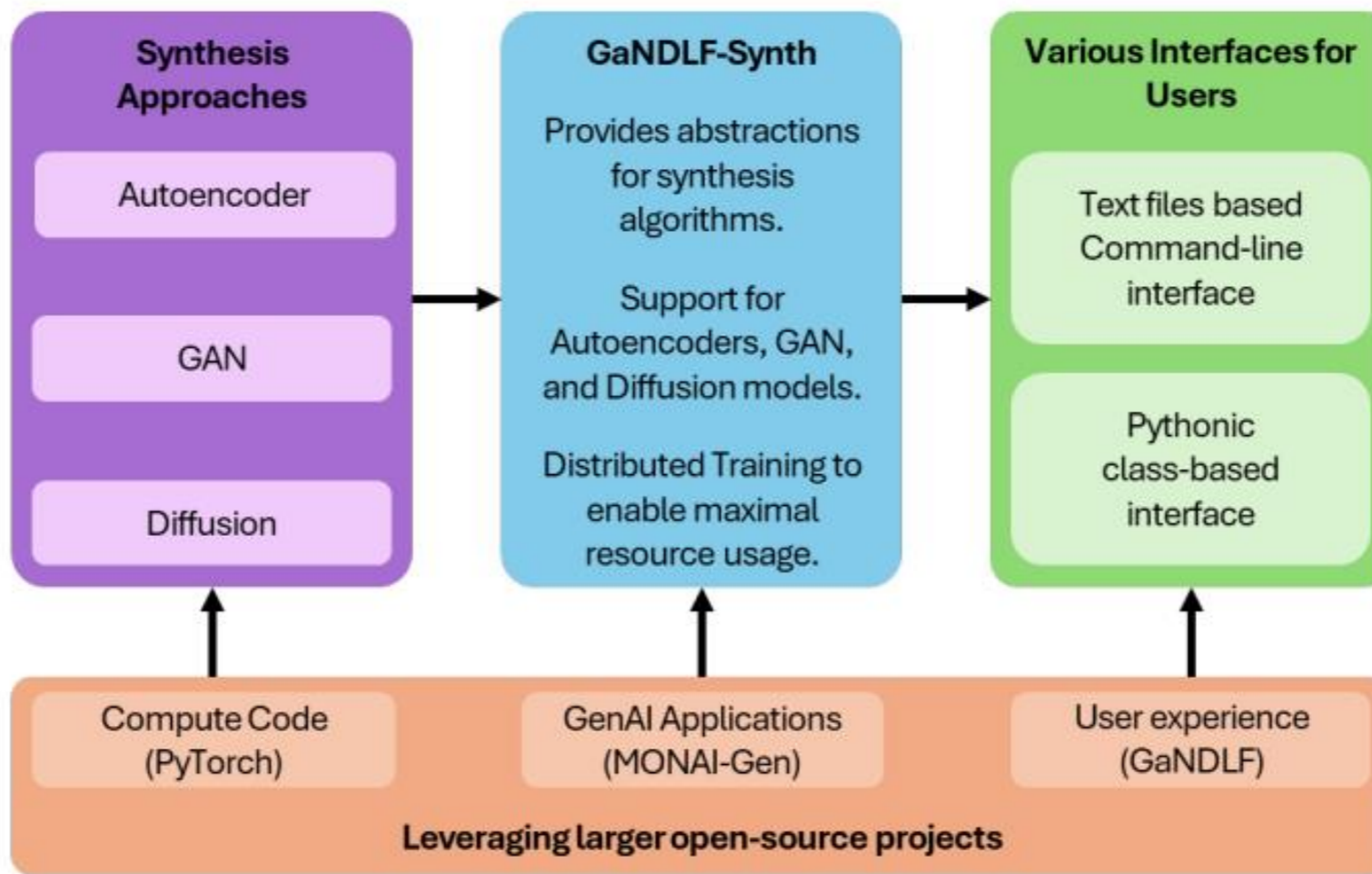
Proposed solution

- **GaNDLF-Synth:** low-code CLI application
- **Simple setup:** based on simple YAML file configuration files
- **End-to-end solution:** data preprocessing, augmentation, data splits, training and inference
- **Scalability - Pytorch Lightning support:** multi-GPU, multinode, DeepSpeed support for running large models
- **Support for multiple neural network architectures:** autoencoders, GANs, diffusion models
- **Extensibility:** modular design, simple to implement custom solutions
- **Validation and robustness:** CI/CD, testing

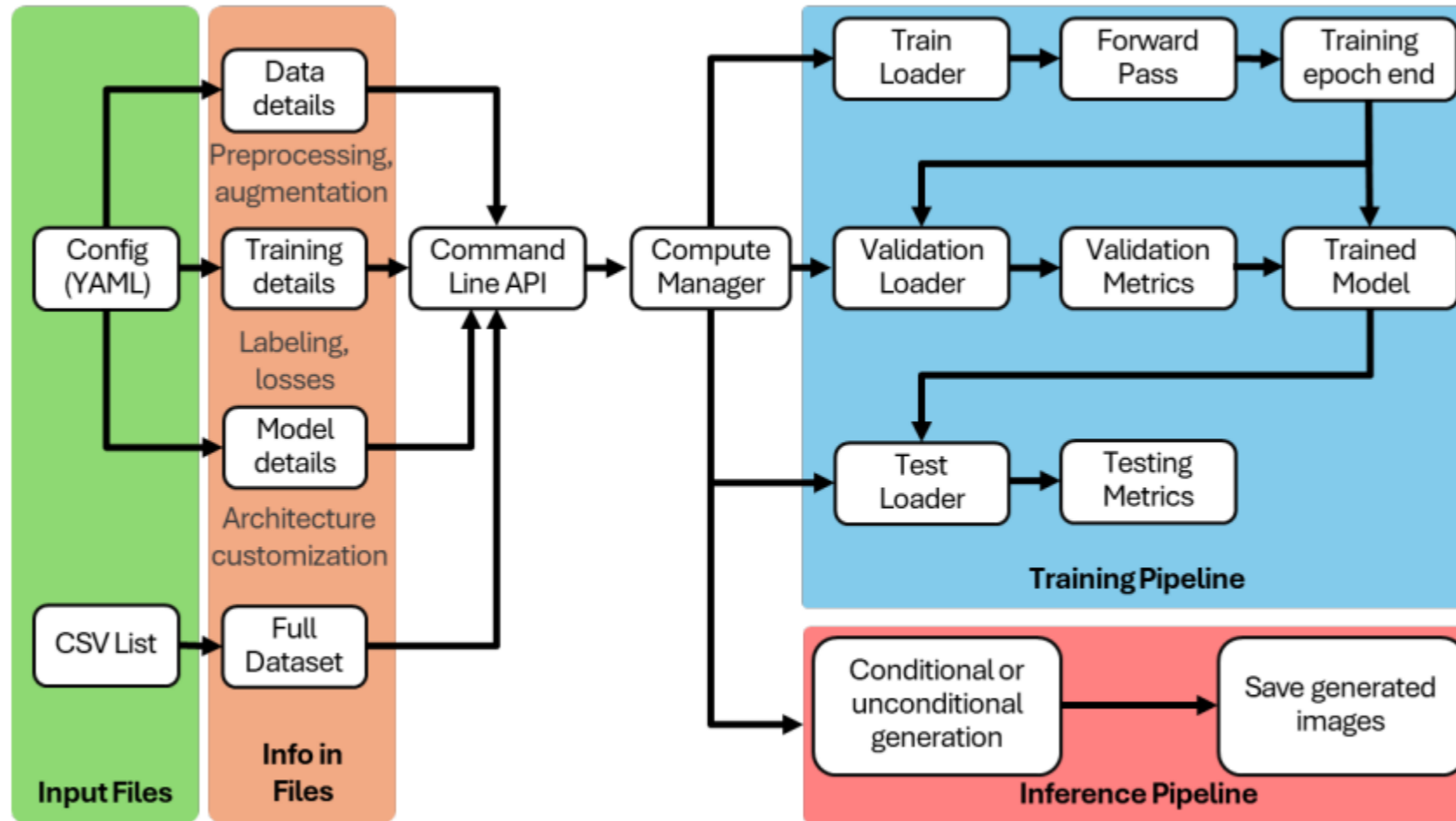
Configuration file example

```
data_augmentation: {}      dataloader_config:      model_config:
data_postprocessing: {}   inference:              architecture:
data_preprocessing:      drop_last: false       num_eval_timesteps: 1
  test:                  num_workers: 0        num_train_timesteps: 1
    resize:              pin_memory: false     converter_type: soft
      - 64                shuffle: false        labeling_paradigm: unlabeled
      - 64                test:                  losses:
train:                   drop_last: false      name: mse
  resize:                num_workers: 0       model_name: ddpm
    - 64                  pin_memory: false    n_channels: 2
    - 64                  shuffle: false       n_dimensions: 2
val:                     train:                norm_type: batch
  resize:                drop_last: false     optimizers:
    - 64                  num_workers: 0      lr: 0.0001
    - 64                  pin_memory: false    name: adam
inference:               validation:           tensor_shape:
  resize:                drop_last: false     - 64
    - 64                  num_workers: 0      - 64
    - 64                  pin_memory: false    schedulers:
                           shuffle: false      type: triangle
                           step_size: 2
```

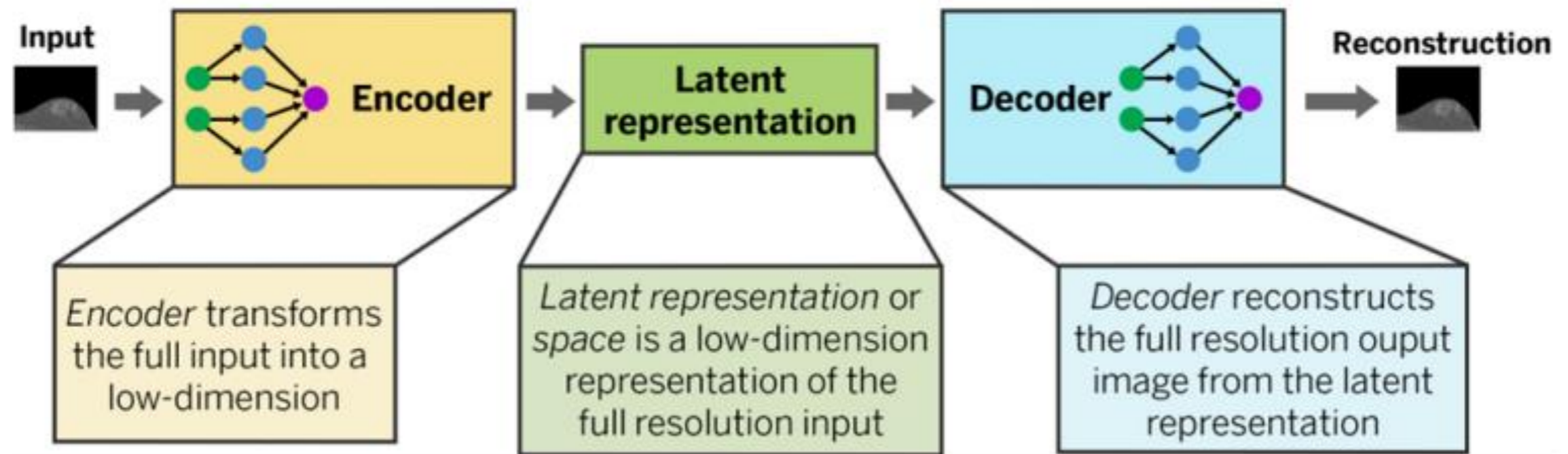

GaNDLF-Synth core principles



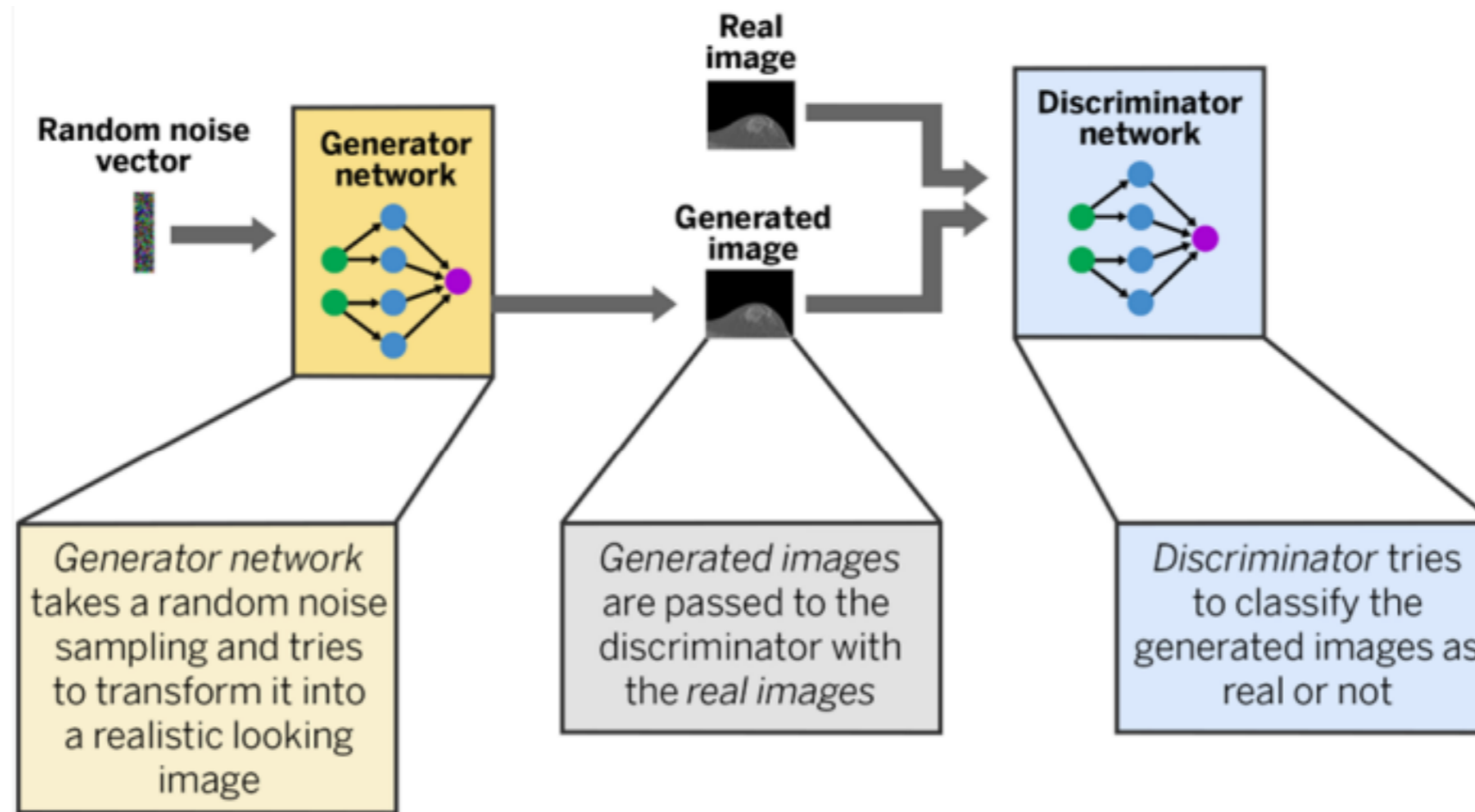
Application architecture



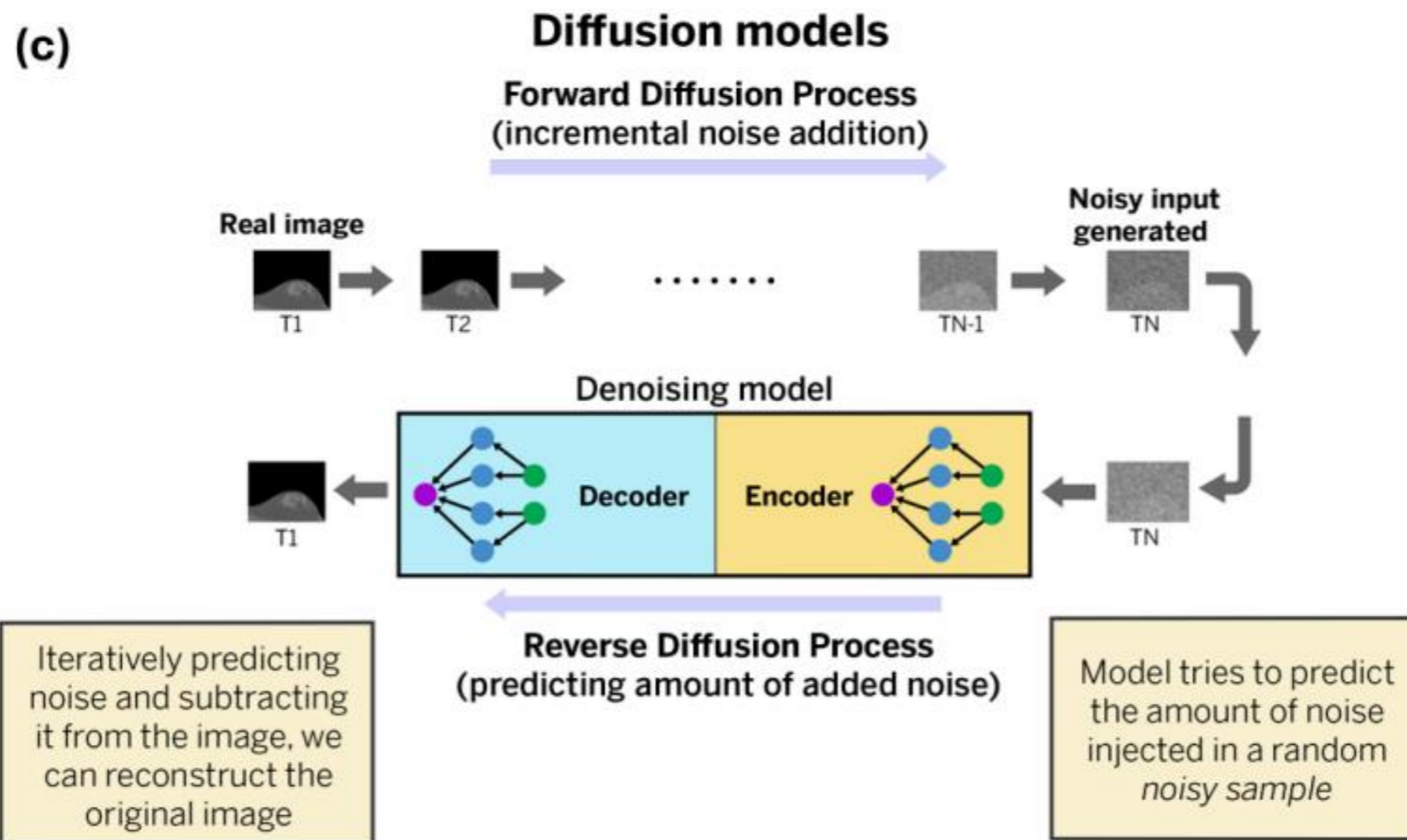
Supported architectures - autoencoders



Supported architectures - GANs



Supported architectures – diffusion models



Conclusions and future work

- Integration of new model architectures
- Federated learning support
- Maintenance and refactoring
- github.com/mlcommons/GaNDLF-Synth



Thank you!