

# Accelerating Science with the Generative Toolkit for Scientific Discovery (GT4SD)

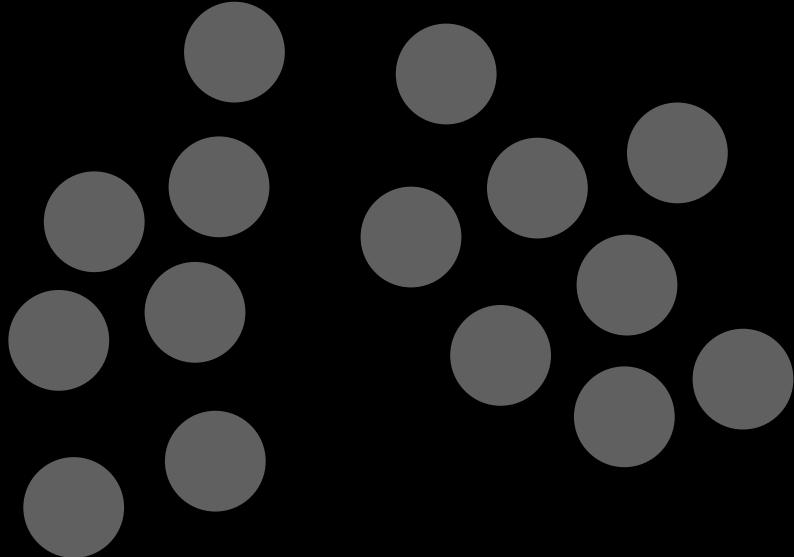


[GT4SD/gt4sd-core](https://github.com/GT4SD/gt4sd-core)

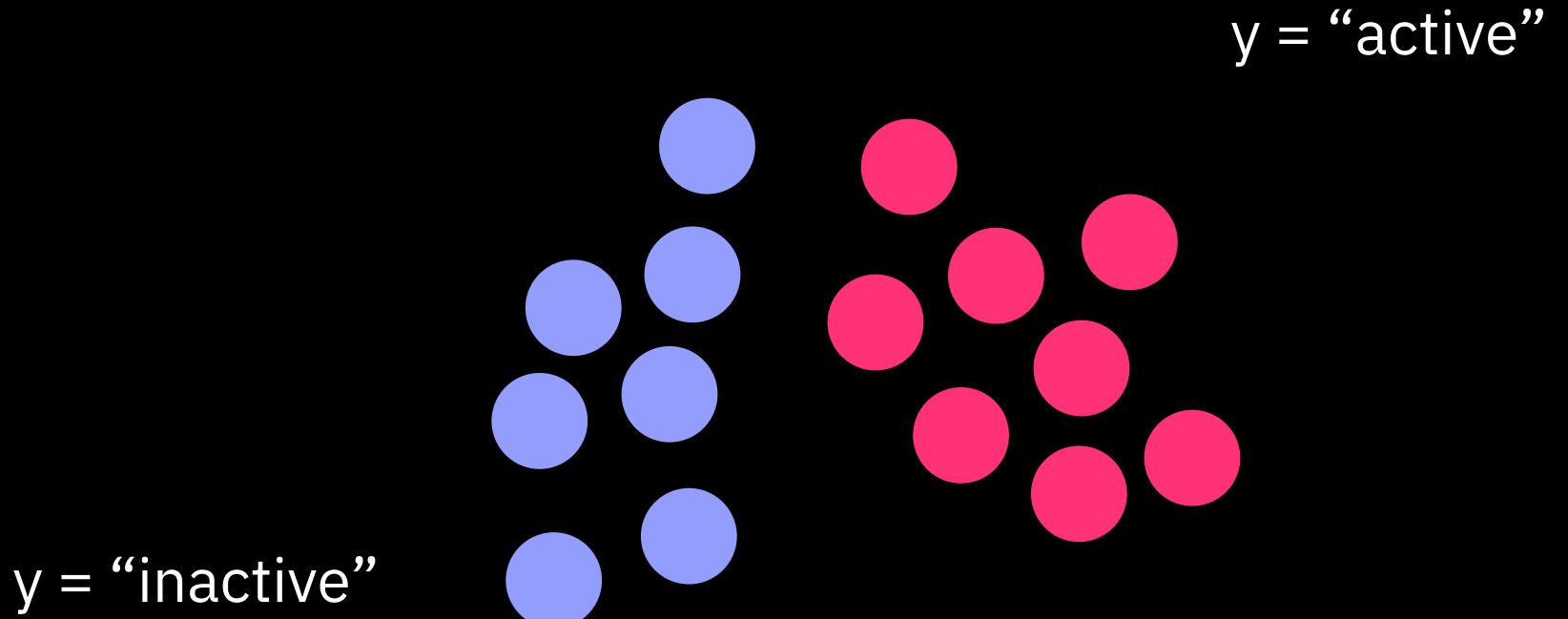


IBM Research

# Introduction to generative models

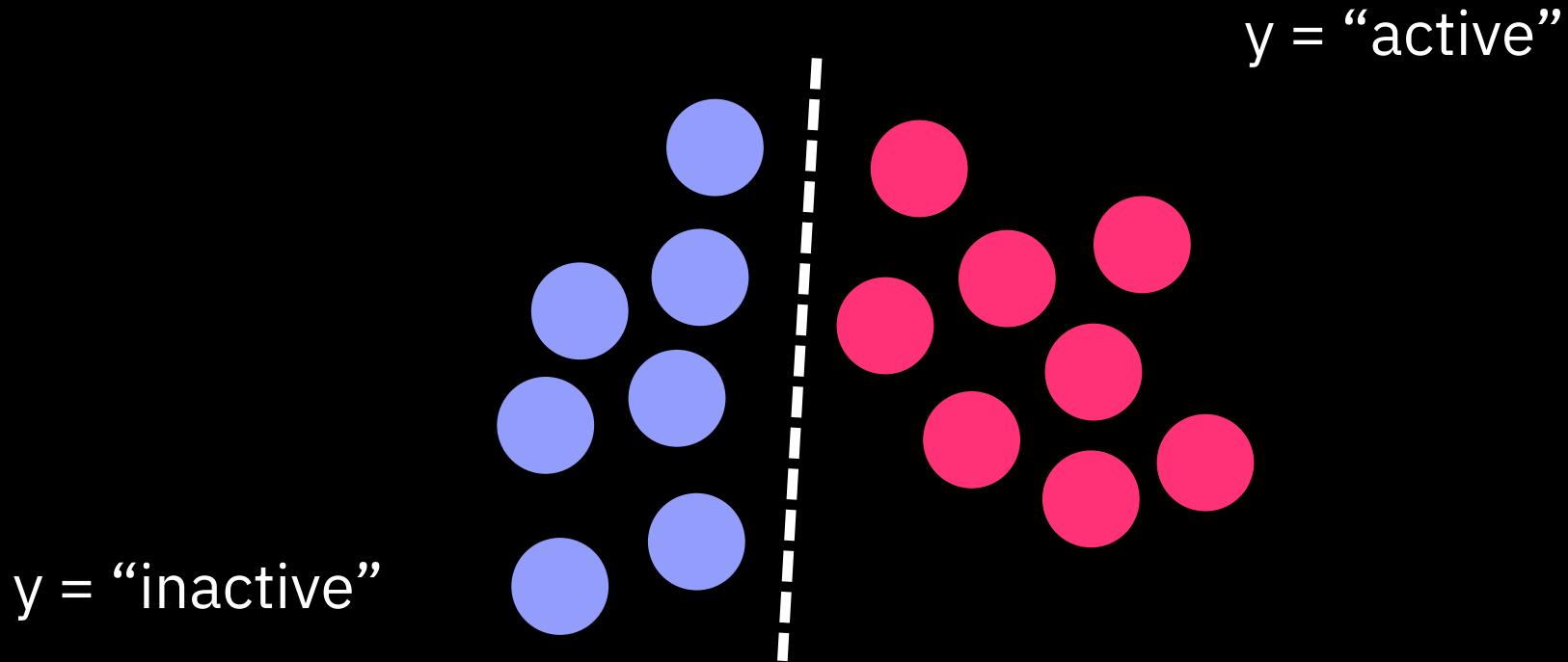


# Introduction to generative models



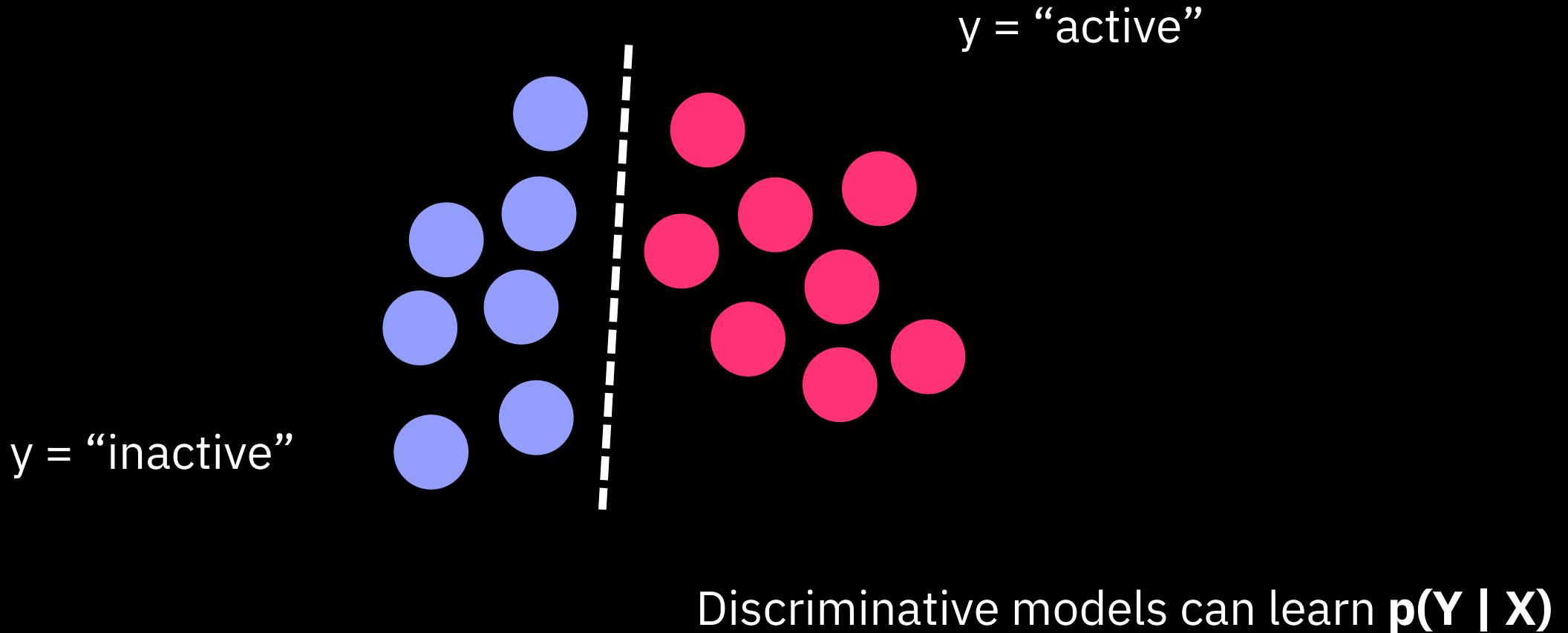
Considering data points (**X**) and  
associated labels (**Y**)

# Introduction to generative models

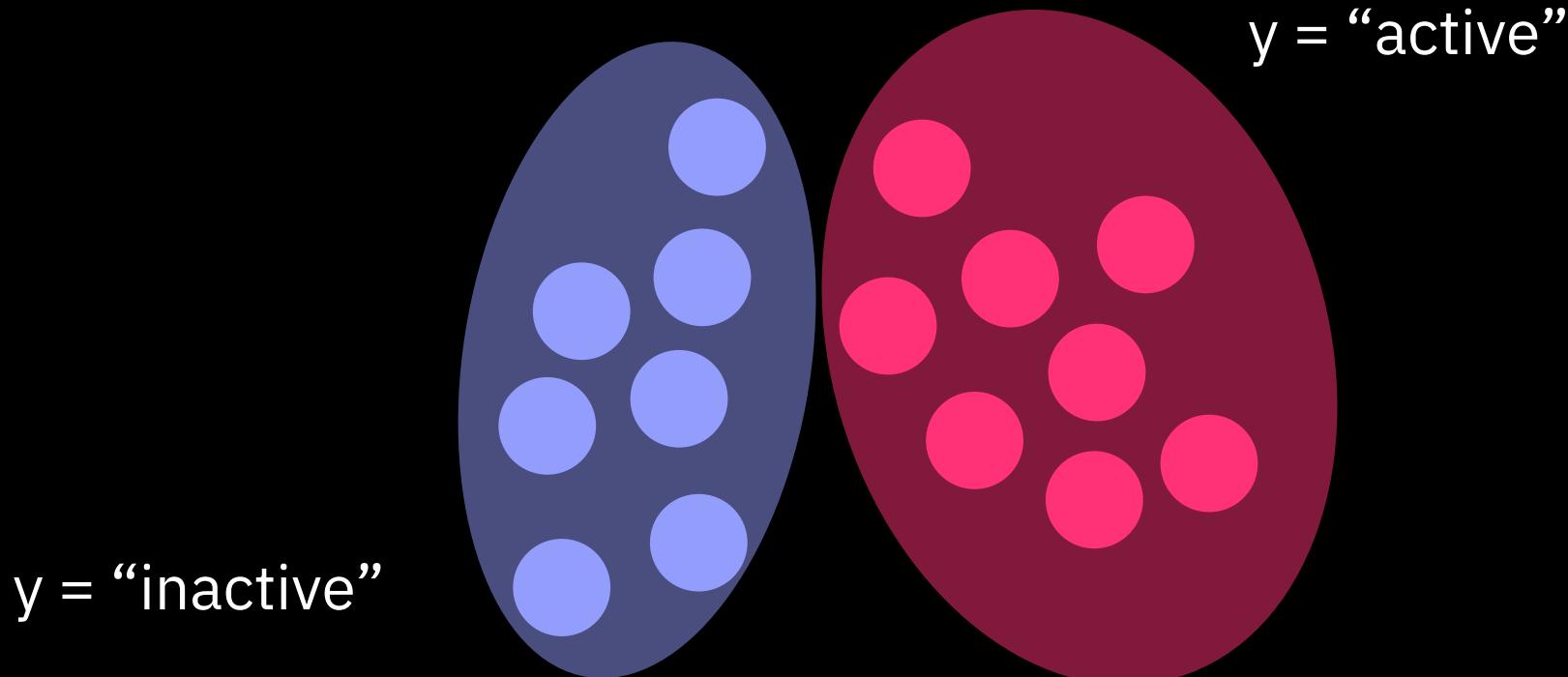


Discriminative models can learn  
boundaries

# Introduction to generative models

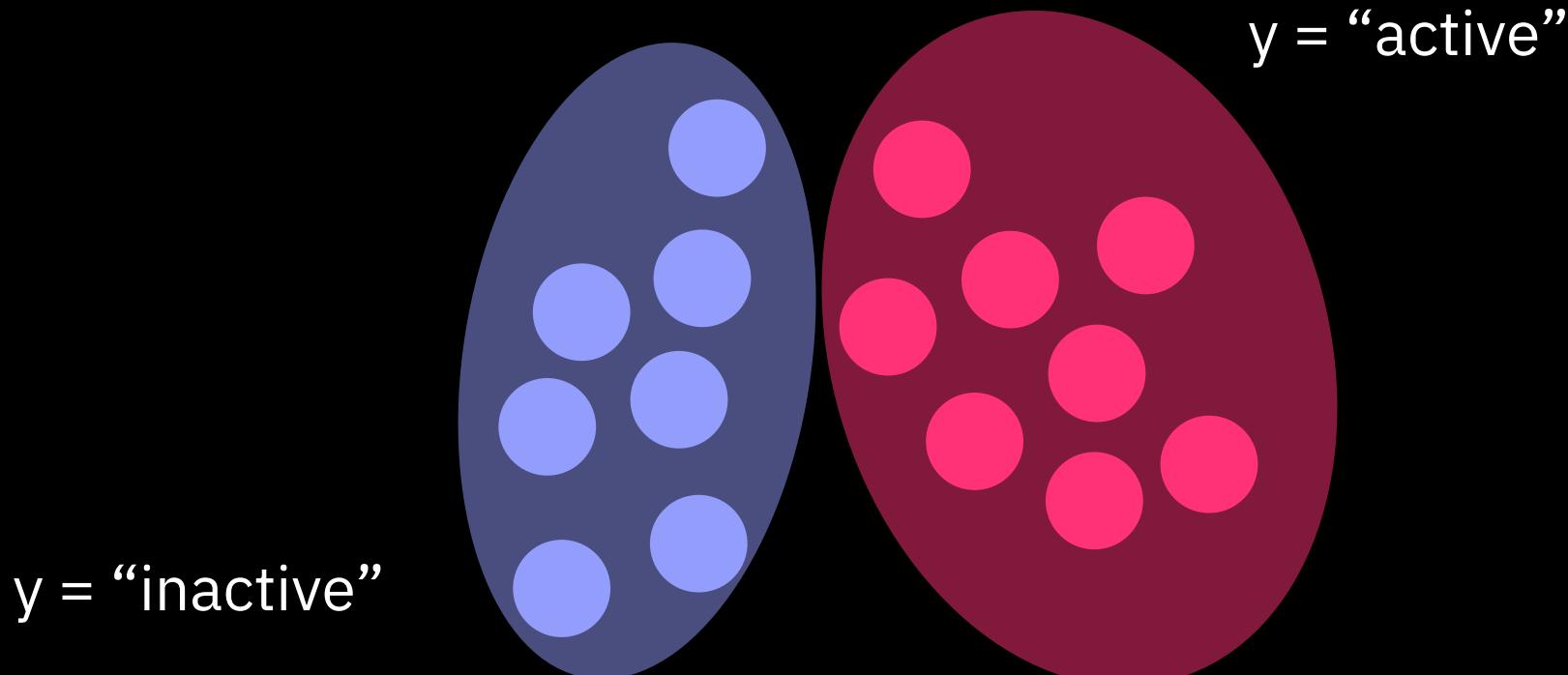


# Introduction to generative models



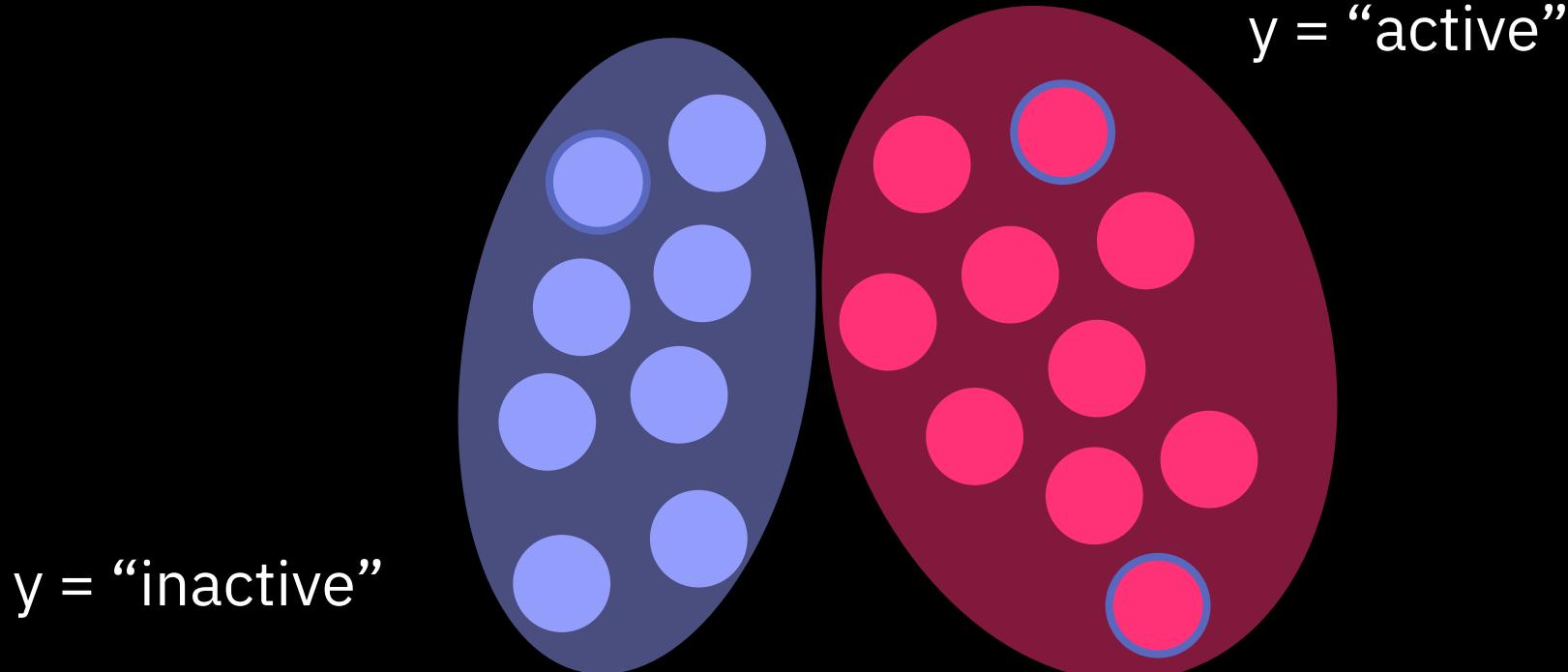
Generative can learn joint data-labels distributions

# Introduction to generative models



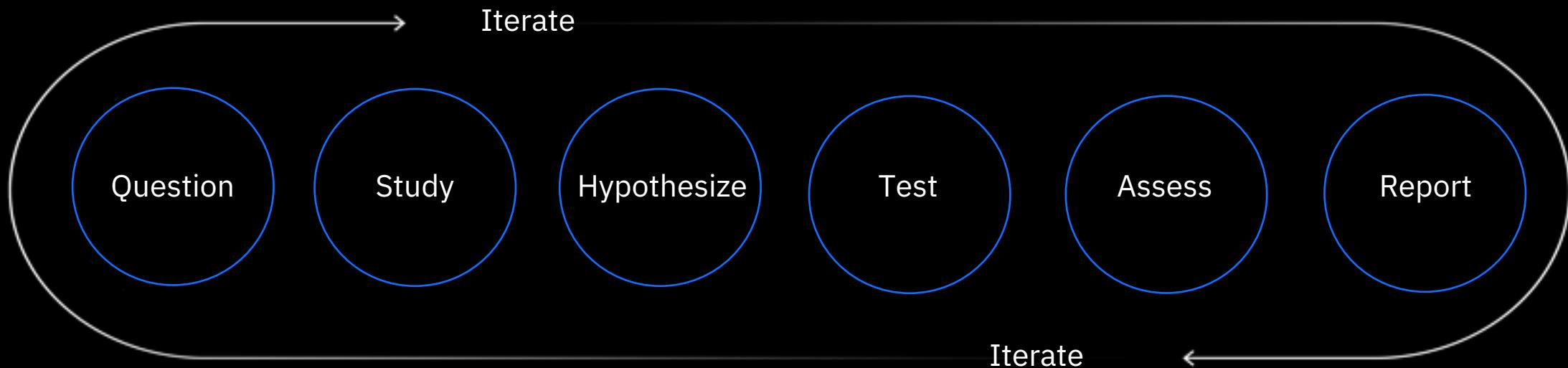
Generative models learn  $p(\mathbf{X}, \mathbf{Y})$

# Introduction to generative models

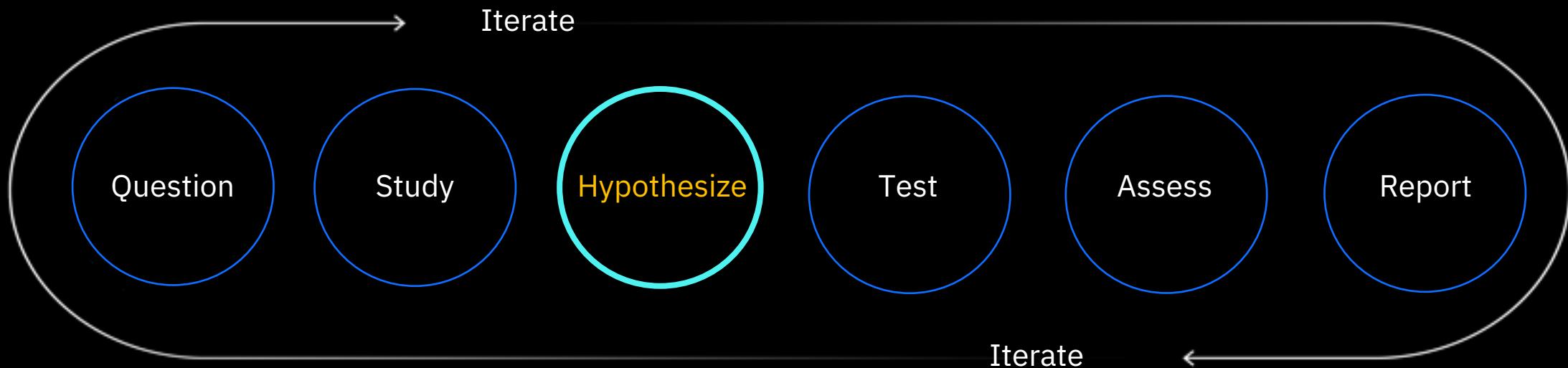


This allows to generate and create novel instances sampling from the underlying distributions

Generative models have the power to supercharge various aspects of the scientific method



We decided to start focusing on the step of accelerating hypothesis generation



# Generative Toolkit for Scientific Discovery (GT4SD)



GT4SD makes generative AI algorithms and models easier to use in scientific discovery



[GT4SD/gt4sd-core](#)

## 1. Train generative models

```
gt4sd-trainer --training_pipeline_name paccmann-vae-trainer --epochs 25
```

## 2. Create inference pipelines

```
gt4sd-saving --training_pipeline_name paccmann-vae-trainer --model_path
```

## 3. Run inference pipelines

```
gt4sd-inference --algorithm_name PaccMannGP --algorithm_application Pacc
```

## 4. Share your models with the community

```
gt4sd-upload --training_pipeline_name paccmann-vae-trainer --model_path
```

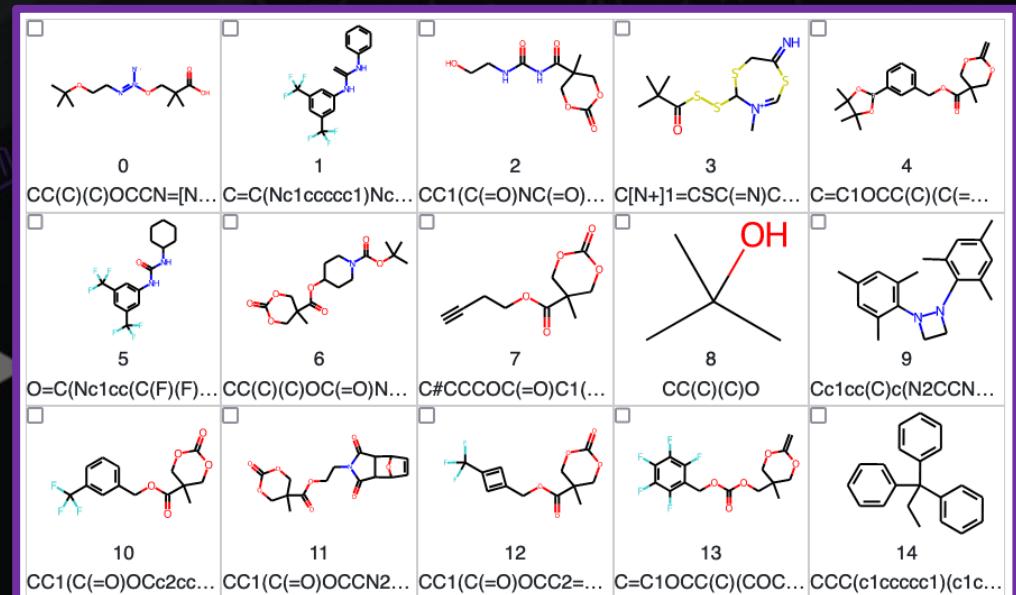
Study

Hypothesize

Test

Applications include hypothesis generation for inverse design and discovery of materials (25+ models are available)

Molecules generated with GT4SD workflows



# Generative Toolkit for Scientific Discovery (GT4SD)



GT4SD makes generative AI algorithms and models easier to use in scientific discovery



[GT4SD/gt4sd-core](#)



```
from gt4sd.algorithms.registry import ApplicationsRegistry
# target definition (can be None)
target = ...
algorithm = ApplicationsRegistry.get_application_instance(
    target=target,
    algorithm_type="algorithm_type",
    domain="materials",
    algorithm_name="AlgorithmName",
    algorithm_application="AlgorithmApplication"
    # include additional configuration parameters as **kwargs
)
# get 50 samples from the algorithm
generated_samples = list(algorithm.sample(50))
```

```
# train a generative algorithm
gt4sd-trainer --training_pipeline_name ${TRAINING_PIPELINE_NAME} \
    --a_parameter 250 --another_parameter 4 --batch_size

# save it for usage
gt4sd-saving --training_pipeline_name ${TRAINING_PIPELINE_NAME} \
    --a_path /path/to/artifacts --target_version v1 \
    --algorithm_application AlgorithmApplication

# generate samples
gt4sd-inference --algorithm_name AlgorithmName \
    --algorithm_application AlgorithmApplication \
    --algorithm_version v1 --number_of_samples 25
```

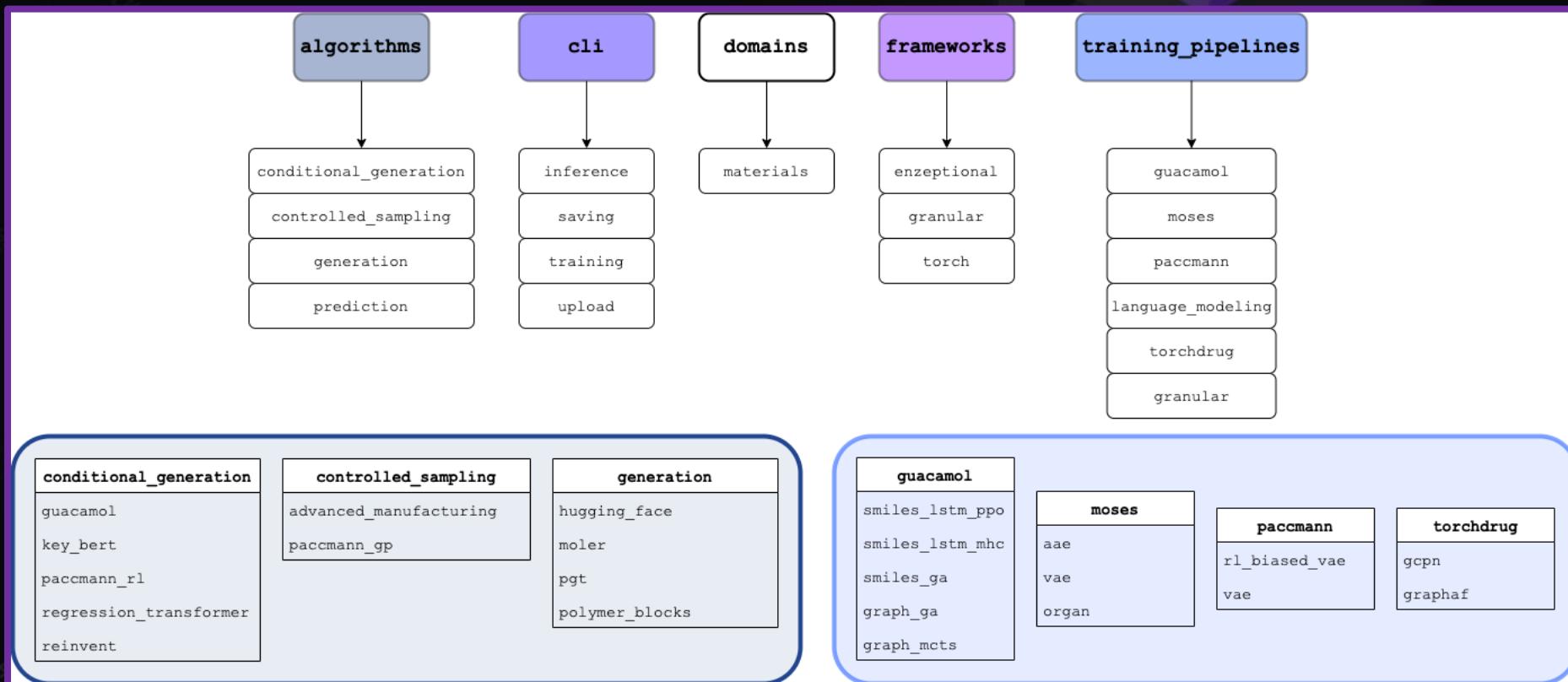
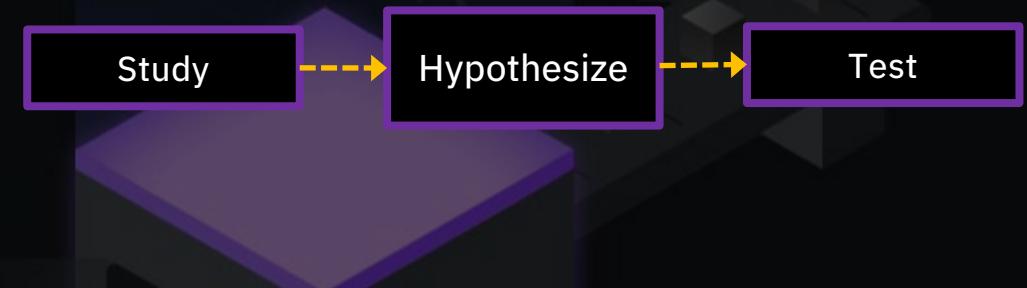
# Generative Toolkit for Scientific Discovery (GT4SD)



GT4SD makes generative AI algorithms and models easier to use in scientific discovery



[GT4SD/gt4sd-core](#)



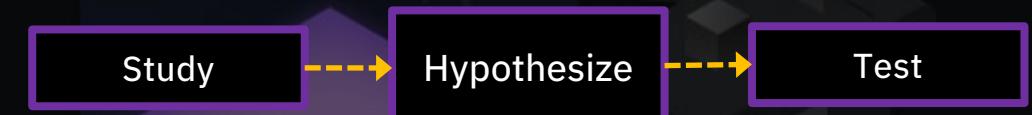
# Generative Toolkit for Scientific Discovery (GT4SD)



GT4SD makes generative AI algorithms and models easier to use in scientific discovery



[GT4SD/gt4sd-core](#)



algorithms    cli    domains    frameworks    training\_pipelines

## Supported packages

Beyond implementing various generative modeling inference and training pipelines GT4SD is designed to provide a high-level API that implement an harmonized interface for several existing packages:

- [GuacaMol](#): inference pipelines for the baselines models and training pipelines for LSTM models.
- [Moses](#): inference pipelines for the baselines models and training pipelines for VAEs and Organ.
- [TorchDrug](#): inference and training pipelines for GCPN and GraphAF models. Training pipelines support custom datasets as well as datasets native in TorchDrug.
- [MoLeR](#): inference pipelines for MoLeR (MOlecle-LEvel Representation) generative models for de-novo and scaffold-based generation.
- [TAPE](#): encoder modules compatible with the protein language models.
- [PaccMann](#): inference pipelines for all algorithms of the PaccMann family as well as training pipelines for the generative VAEs.
- [transformers](#): training and inference pipelines for generative models from [HuggingFace Models](#)

conditional_generation	controlled_generation
guacamol	advanced_generation
key_bert	paccmann
paccmann_rl	
regression_transformer	
reinvent	

# Generative Toolkit for Scientific Discovery (GT4SD)

Textual representation of molecules  
(SMILES/SELFIES) combined w/ GAN/VAE/GA

Study

Hypothesize

Test

frameworks

training\_pipelines

Beyond implementing various generative modeling inference and training pipelines GT4SD is designed to provide a high-level API that implement an harmonized interface for several existing packages:

- [GuacaMol](#): inference pipelines for the baselines models and training pipelines for LSTM models.
- [Moses](#): inference pipelines for the baselines models and training pipelines for VAEs and Organ.
- [TorchDrug](#): inference and training pipelines for GCPN and GraphAF models. Training pipelines support custom datasets as well as datasets native in TorchDrug.
- [MoLeR](#): inference pipelines for MoLeR (MOlecle-LEvel Representation) generative models for de-novo and scaffold-based generation.
- [TAPE](#): encoder modules compatible with the protein language models.
- [PaccMann](#): inference pipelines for all algorithms of the PaccMann family as well as training pipelines for the generative VAEs.
- [transformers](#): training and inference pipelines for generative models from [HuggingFace Models](#)

conditional\_generation

- guacamol
- key\_bert
- paccmann\_rl
- regression\_transformer
- reinvent

controlled

- advanced
- paccmann

general

- conditional
- controlled
- general
- predicted

# Generative Toolkit for Scientific Discovery (GT4SD)

Graph representation of molecules combined w/  
GCN/VAE/GA

Study → Hypothesize → Test

frameworks training\_pipelines

Beyond implementing various generative modeling inference and training pipelines GT4SD is designed to provide a high-level API that implement an harmonized interface for several existing packages:

- [GuacaMol](#): inference pipelines for the baselines models and training pipelines for LSTM models.
- [Moses](#): inference pipelines for the baselines models and training pipelines for VAEs and Organ.
- [TorchDrug](#): inference and training pipelines for GCPN and GraphAF models. Training pipelines support custom datasets as well as datasets native in TorchDrug.
- [MoLeR](#): inference pipelines for MoLeR (MOlecle-LEvel Representation) generative models for de-novo and scaffold-based generation.
- [TAPE](#): encoder modules compatible with the protein language models.
- [PaccMann](#): inference pipelines for all algorithms of the PaccMann family as well as training pipelines for the generative VAEs.
- [transformers](#): training and inference pipelines for generative models from [HuggingFace Models](#)

conditional\_generation  
guacamol  
key\_bert  
paccmann\_rl  
regression\_transformer  
reinvent

controlled  
advanced  
paccmann

# GT4SD enables complete discovery workflows

## Context:

This molecule was proposed by GENTRL (a deep generative model) as DDR1-inhibitor. It showed favorable pharmacokinetics in mice.

For details see Zhavoronkov et al. (2019)

QED: 0.38

## Task:

Find a **similar** molecule (measured by Tanimoto similarity) with an improved **drug-likeness** (QED) score.

### Step 1: Investigate the chemical space of molecular generative models

#### Graph models

- GCPN
- Graph AF
- MoLeR

#### Language models

C1(C=O)cc(OC)c(O)cc1

- VAE
- AAE
- ORGAN
- GuacaMol & MOSES

### Step 2: Investigate **conditional** generative models that can be primed with **properties** or **substructures** (scaffolds)

#### Models

- REINVENT
- MoLeR
- Regression Transformer

Tanimoto similarity

QED

Desired region

Algorithm: ReINVENT, MoLeR, MoLeR Scaff, Regression/Transformer

MoLeR

Tanimoto: 0.59 QED: 0.47

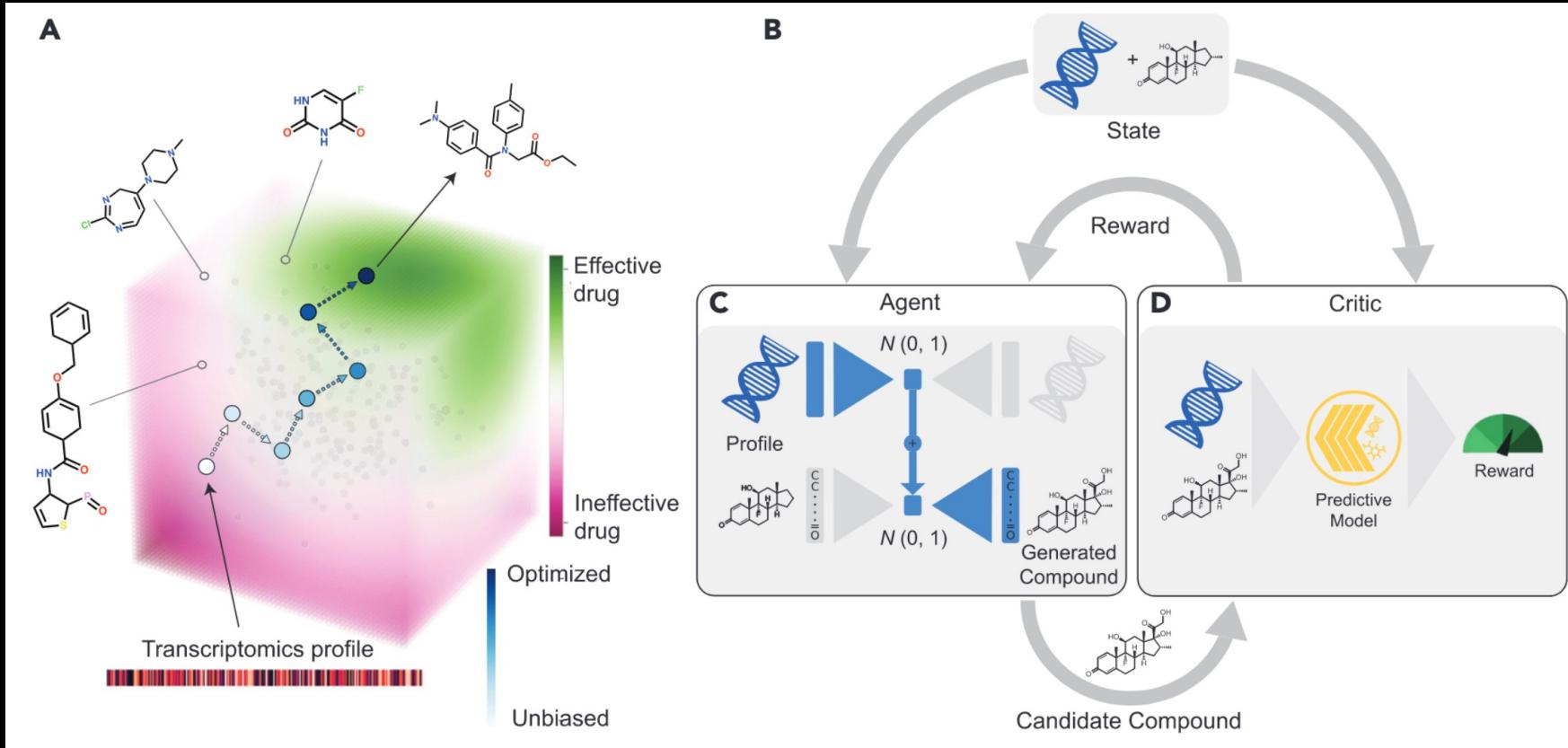
Regression Transformer

Tanimoto: 0.82 QED: 0.39

REINVENT

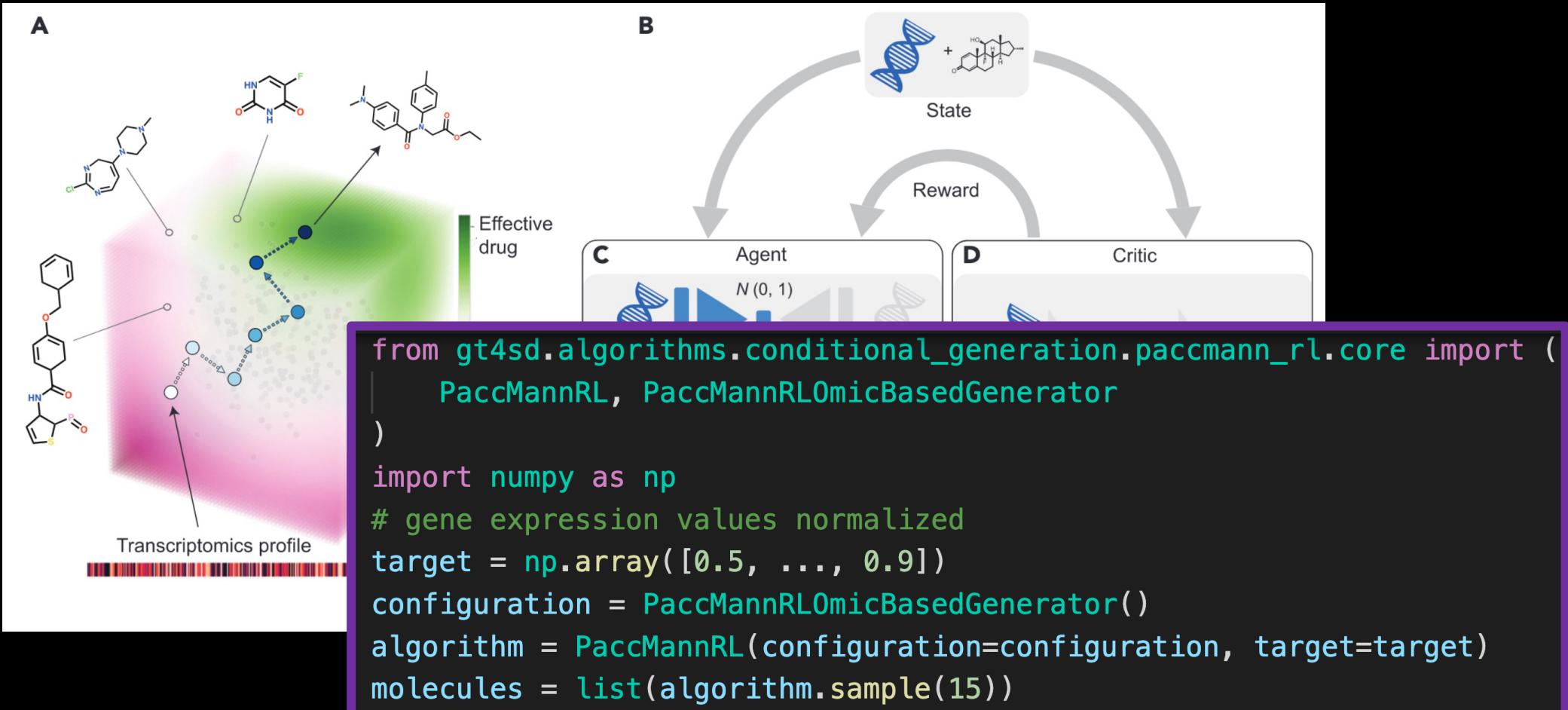
Tanimoto: 0.12 QED: 0.72

# String-based (1D) Generative Molecule Creation - Gene Expression-driven Hit-like Molecule Design



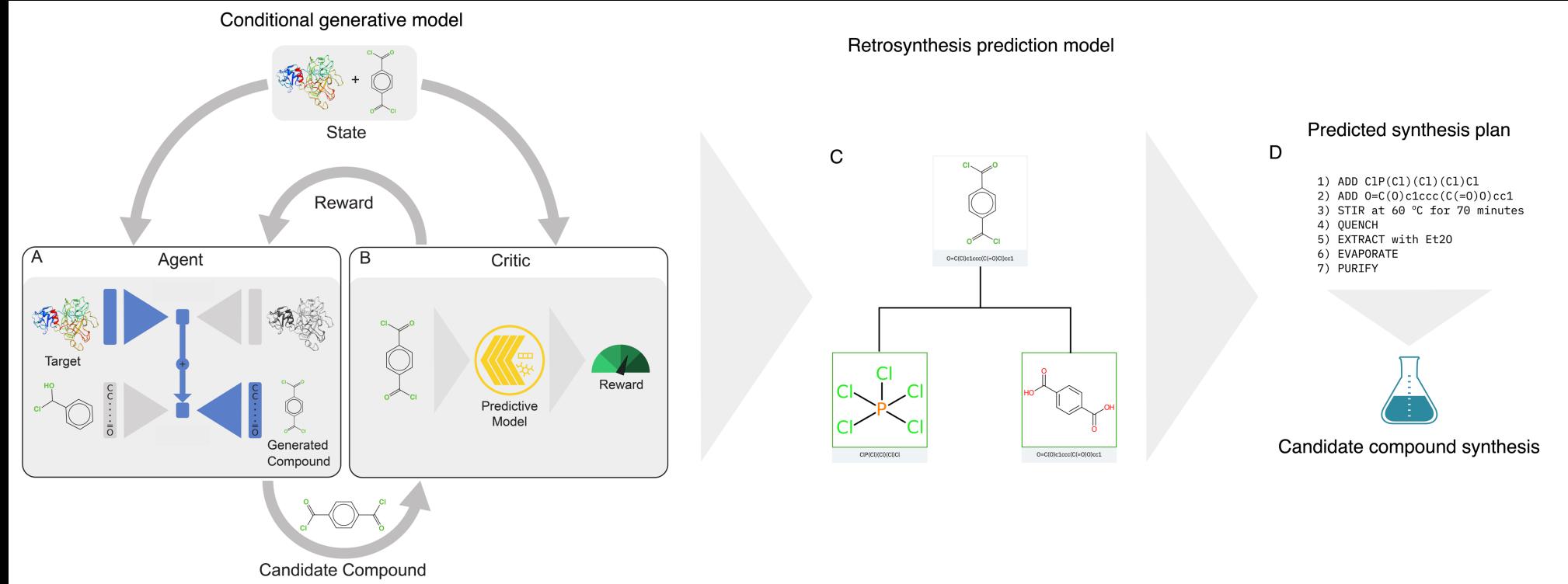
Born *et al.* (2021). PaccMannRL: De novo generation of hit-like anticancer molecules from transcriptomic data via reinforcement learning. *iScience*, 24, 4.

# String-based (1D) Generative Molecule Creation - Gene Expression-driven Hit-like Molecule Design



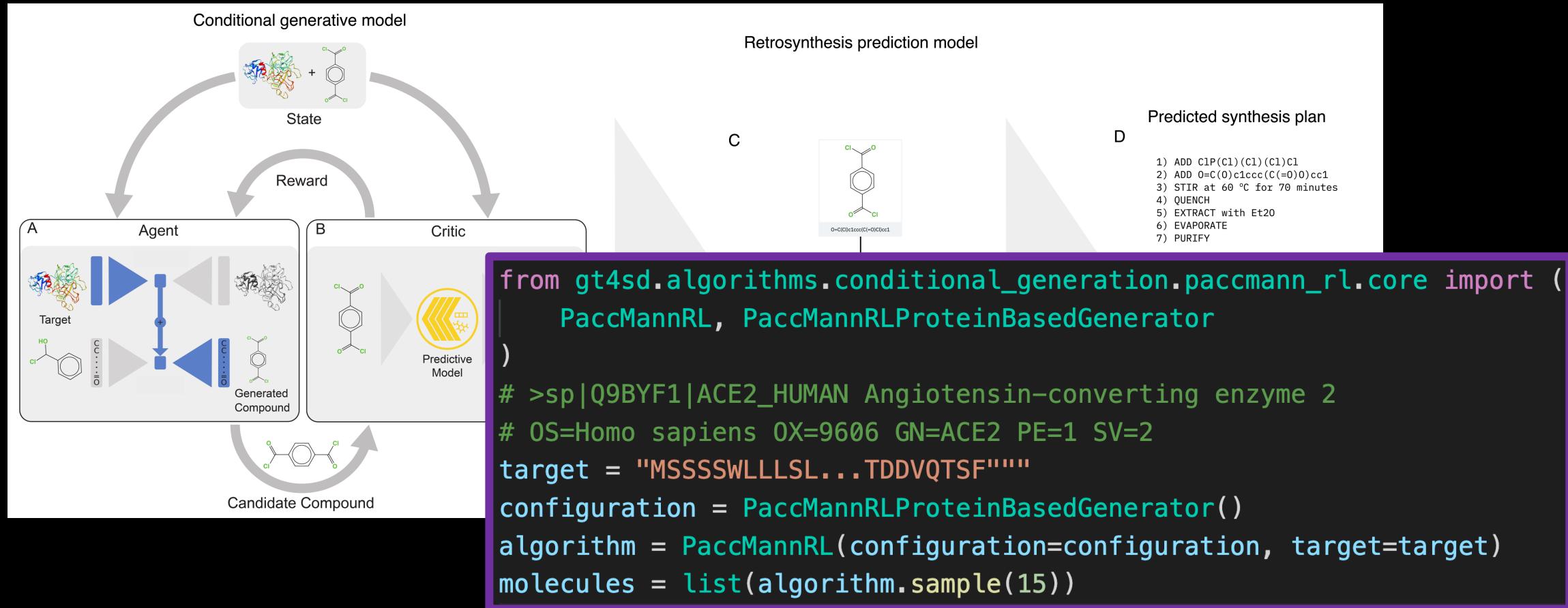
Born *et al.* (2021). PaccMannRL: De novo generation of hit-like anticancer molecules from transcriptomic data via reinforcement learning. *iScience*, 24, 4.

# String-based (1D) Generative Molecule Creation - Data-driven Ligand Design and Synthesis



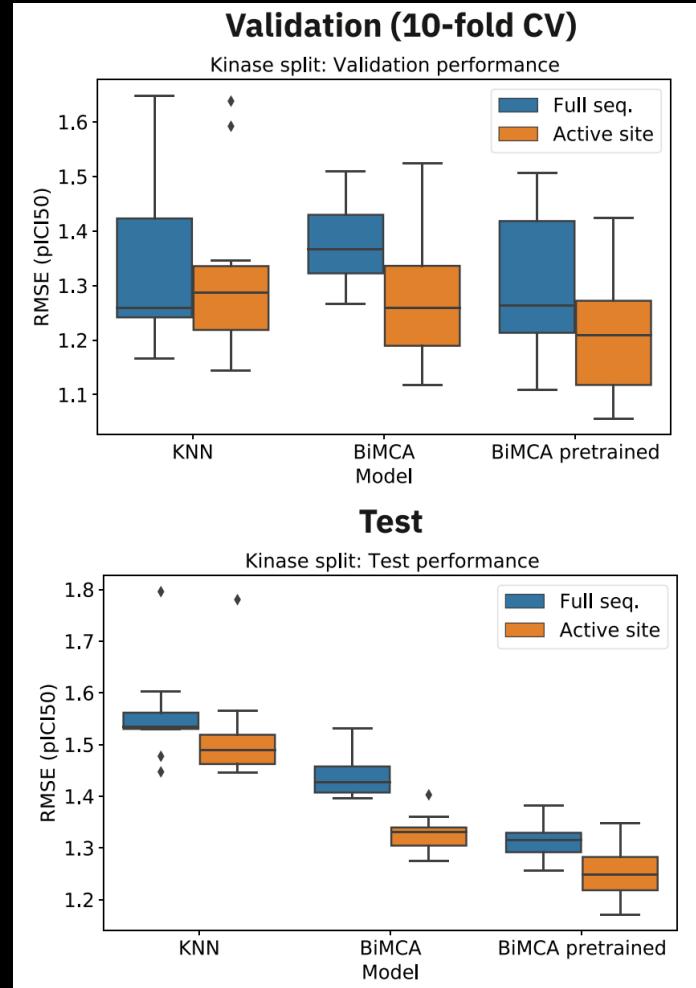
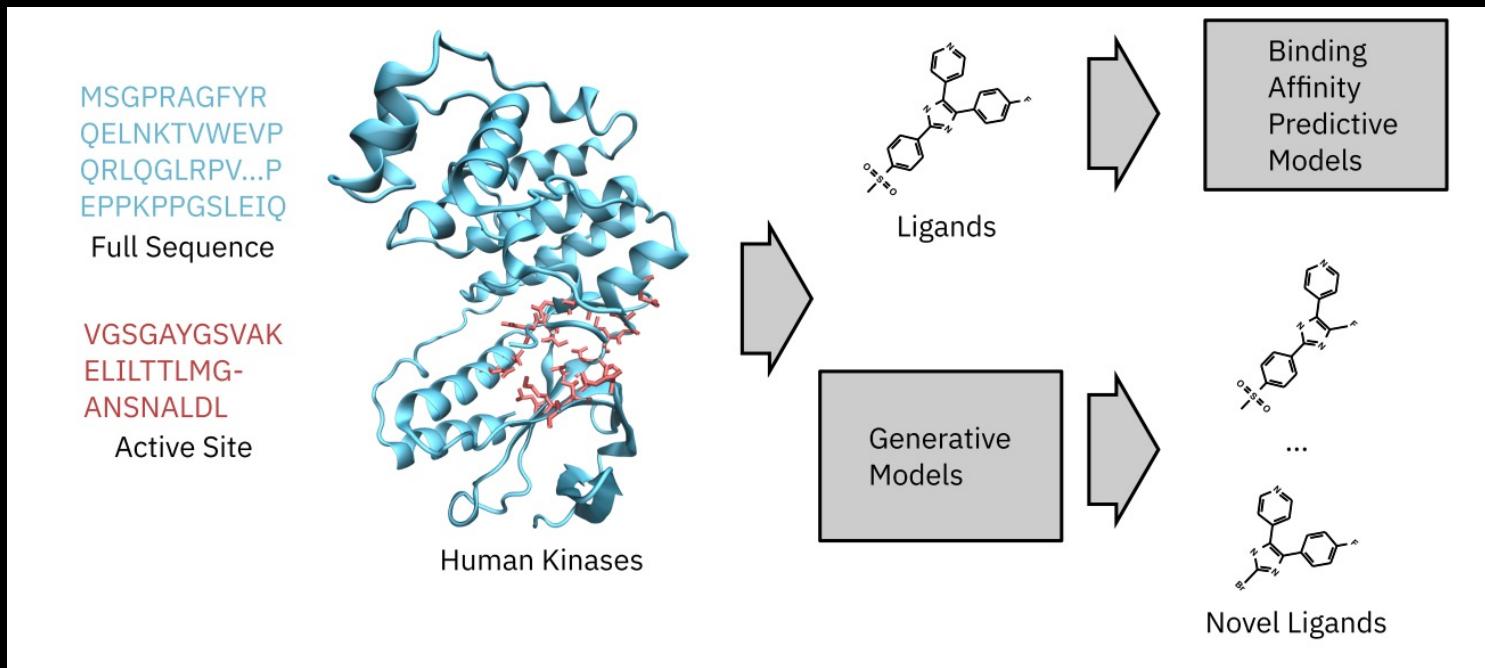
Born *et al.* (2021). Data-driven Molecular Design for Discovery and Synthesis of Novel Ligands-A case study on SARS-CoV-2. *MLST*, 2, 2.

# String-based (1D) Generative Molecule Creation - Data-driven Ligand Design and Synthesis



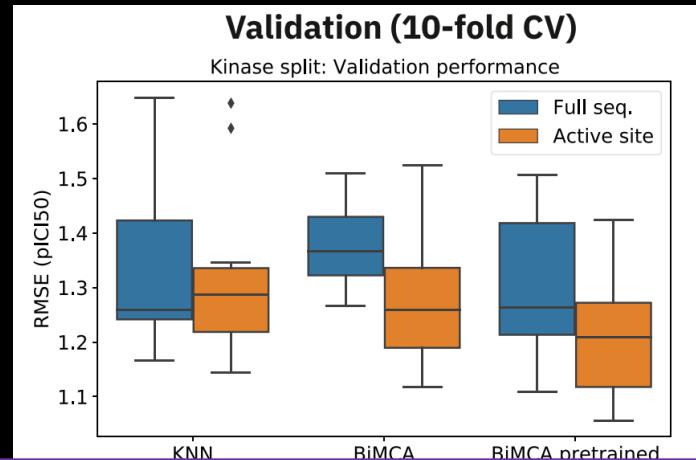
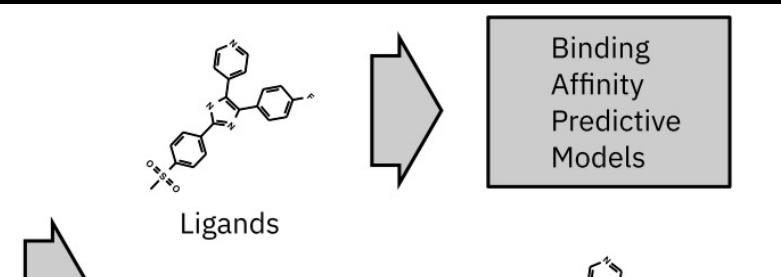
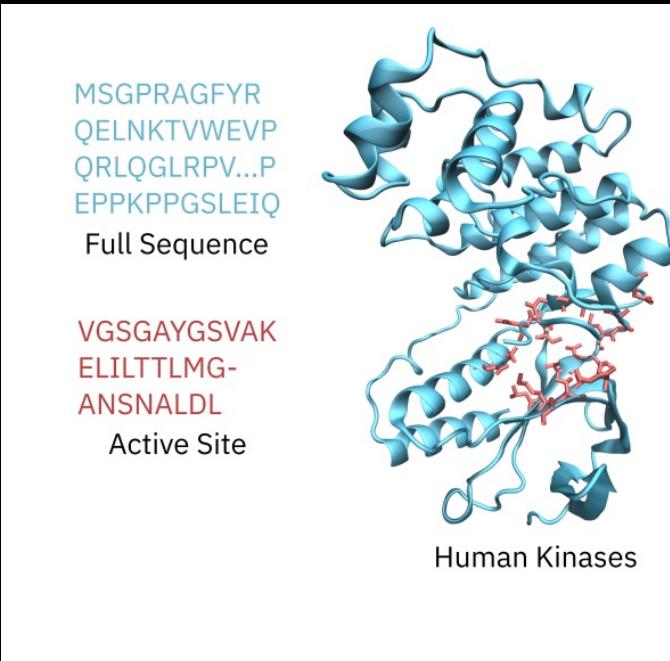
Born *et al.* (2021). Data-driven Molecular Design for Discovery and Synthesis of Novel Ligands-A case study on SARS-CoV-2. *MLST*, 2, 2.

# String-based (1D) Generative Molecule Creation - 3D Effects by Using Only Active Site Residues



Born *et al.* (2022). Active Site Sequence Representations of Human Kinases Outperform Full Sequence Representations for Affinity Prediction and Inhibitor Generation: 3D Effects in a 1D Model. *JCIM*, 62, 240-257.

# String-based (1D) Generative Molecule Creation - 3D Effects by Using Only Active Site Residues

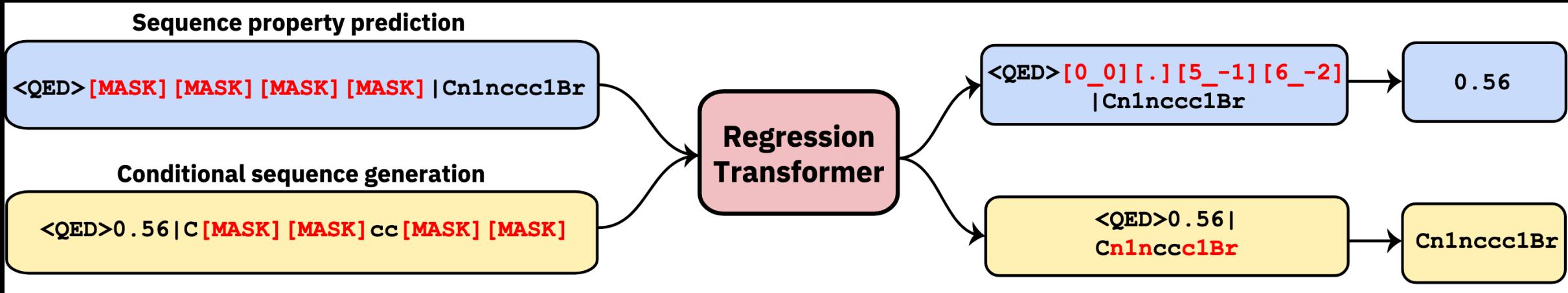


```
from gt4sd.algorithms.controlled_sampling.paccmann_gp.core import (
    PaccMannGP, PaccMannGPGenerator
)

# maximizing drug likeness and synthesizability
target = {"qed": {"weight": 1.0}, "sa": {"weight": 1.0}}
configuration = PaccMannGPGenerator()
algorithm = PaccMannGP(configuration=configuration, target=target)
molecules = list(algorithm.sample(15))
```

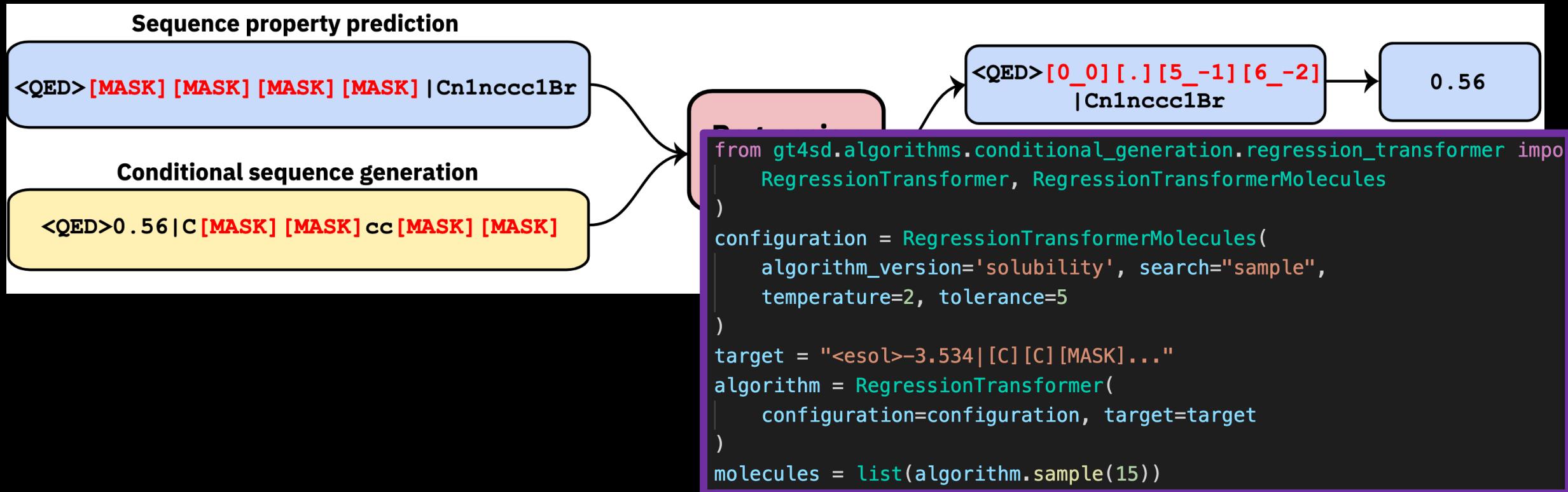
Born *et al.* (2022). Active Site Sequence Representations of Human Kinases Outperform Full Sequence Representations for Affinity Prediction and Inhibitor Generation: 3D Effects in a 1D Model. *JCIM*, 62, 240-257.

# String-based (1D) Generative Molecule Creation - Multi-task Modeling for Conditional Generative Design



Born and Manica (2022). Regression Transformer: Concurrent Conditional Generation and Regression by Blending Numerical and Textual Tokens. *ICLR Workshop Machine Learning for Drug Discovery*

# String-based (1D) Generative Molecule Creation - Multi-task Modeling for Conditional Generative Design



Born and Manica (2022). Regression Transformer: Concurrent Conditional Generation and Regression by Blending Numerical and Textual Tokens. *ICLR Workshop Machine Learning for Drug Discovery*

# Thanks for your attention

## GT4SD (Generative Toolkit for Scientific Discovery)

pypi package 0.44.0 Running tests: style, pytests and entry-points passing License MIT code style black contributions welcome

website live downloads 10k downloads/month 3k launch binder DOI 10.5281/zenodo.6798761

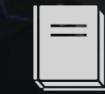
Award 2022 IEEE Open Software Services Award



If you want to know more check out the available resources:



[GT4SD/gt4sd-core](#)



[Pre-print:2207.03928](#)



[Documentation](#)

Contributions are welcome:

[CONTRIBUTING.md](#)

[CODE\\_OF\\_CONDUCT.md](#)