# Machine Learning for *de novo* Molecular Design
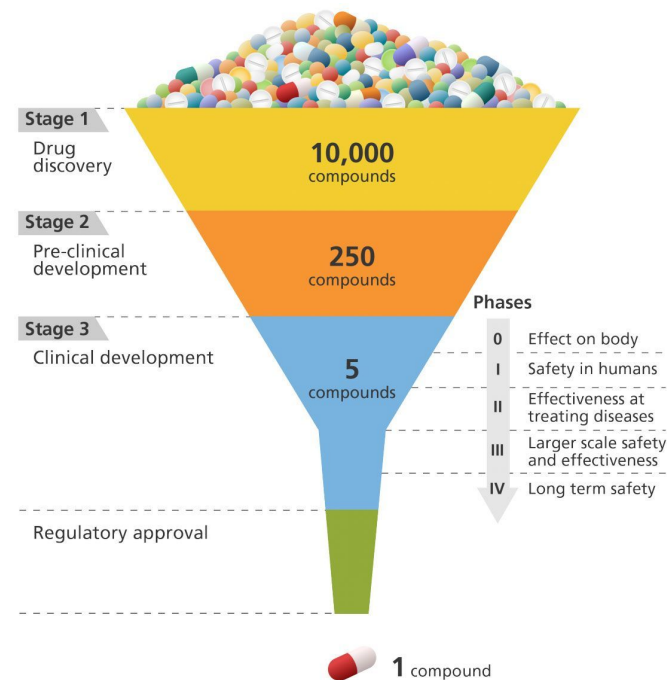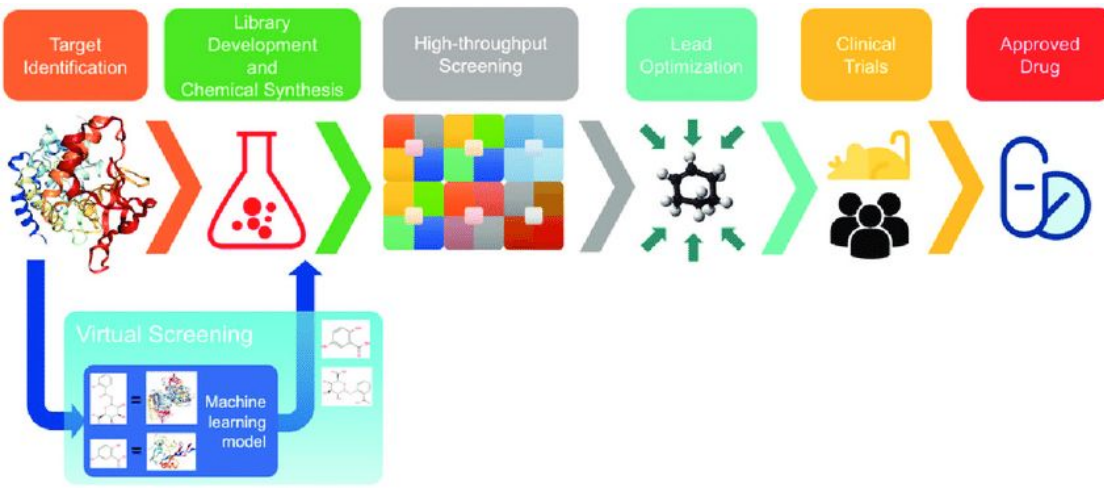
Roi Naveiro

LifeHub

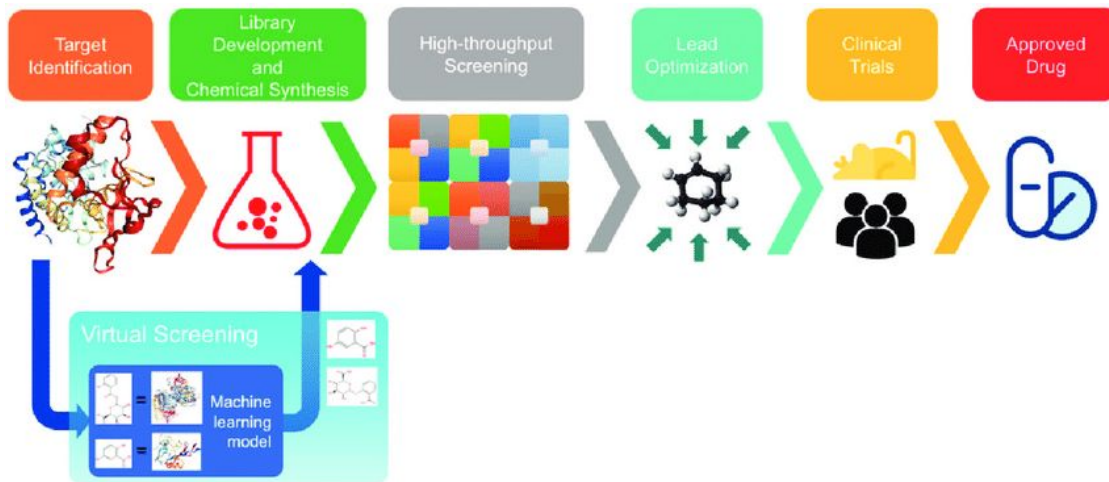My knowledge in chemistry is very (very) basic...

# Why?

# The process of discovering new molecules

- Pharma: average time discovery – market, 13 years

- Outside pharma: 25 years

- Crucial 1st step: **generate pool of candidates**

- Daunting task (e.g. $10^{23}$ – $10^{60}$ drug–like molecules)

# The old way and the soon-to-be-old way

- Old way

  - Human experts propose, synthesize and test (*in vitro*)

- Soon-to-be-old way: high throughput virtual screening (HTVS)

  - Predict properties through computational chemistry...
  - ...leverage rapid **ML-based property predictions**

# De novo molecular design

- Just existing molecules are explored

- Much time lost evaluating bad leads

- Traverse chemical space more "effectively": reach **optimal molecules** with **less evaluations** than brute-force screening

*"De novo molecular design is the process of automatically proposing novel chemical structures that optimally satisfy desired properties"*

**Combinatorial, black-box, stochastic, multi-objective optimization with black-box constraints**

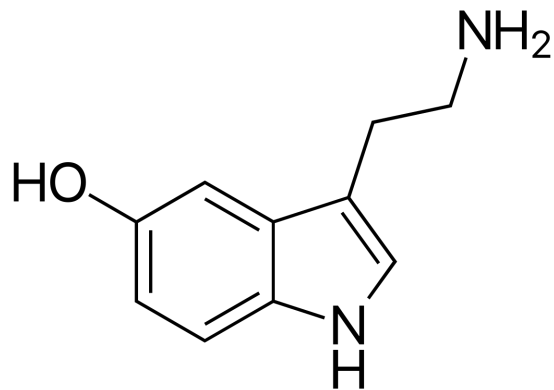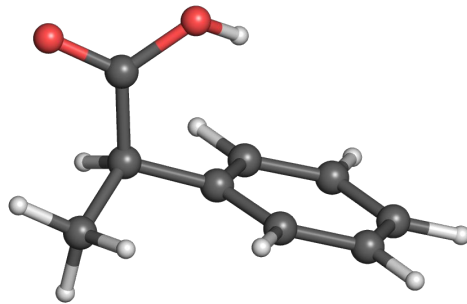# Automatically proposing novel chemical structures

Two main ingredients

- Molecule representation

- Generative model

# Representing molecules

Molecules are **3D QM objects** with: nuclei with defined positions surrounded by electrons described by complex wave-functions

- Digital encoding that serves as input to model

- **Uniqueness and invertibility**

- Trade-off: information lost vs complexity

  - 3D coord. representation (symmetries?)

  - More compact 2D (graph) representation
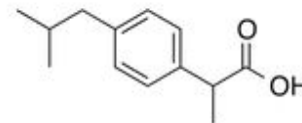
- 1D, 2D and 3D

# 1D representations – SMILEs

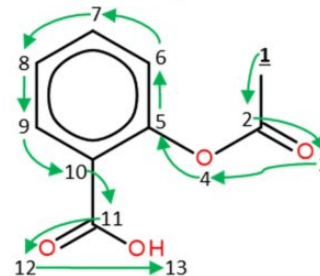**Simplified Molecular Input Line Entry System**

Molecule as graph (bond length and conformational info is lost)

- Graph traversal
- Sequence of ASCII characters
- Non-unique → Canonical SMILES
- One-Hot-Encoding

- Leverage NLP techniques

- SMILE-based methods struggle to generate **valid** molecules
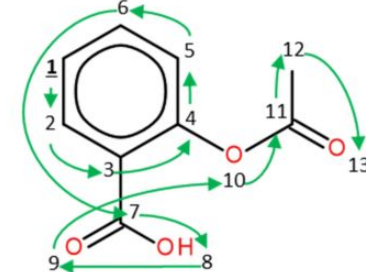- Valid = valency rules
- Learn spurious grammar rules



Ibuprofen

CC(C)Cc1ccc(cc1)C(C)C(O)=O



**a** Canonical representation  **b** Randomized representation

CC(=O)Oc1ccccc1C(=O)O

c1cc(c(cc1)C(O)=O)OC(C)=O

# 2D representations

- Nodes represent atoms
- Edges represent bonds
- Nodes/Edges have associated features (atom number, bond type, etc.)

- Capture connectivity!
- Symmetry invariant representation

- More difficult to generate than sequences
- Taylored algorithms that work with graphs
  (composing transformations on graphs, symmetries?)

- Graph Neural Nets!



Acetaminophen

# 3D representations

- 3D point clouds

$$\mathcal{M} = \{x_i, r_i\}_{i=1}^p$$ where $x_i$ are features and $r_i$ are coordinates.

- Minimal information lost (conformational preferences, bond lengths, etc.)

- Symmetries?
- Too many degrees of freedom

- Generation: sequentially choose pair of atoms, relative position, bond length and angles

# How to generate molecules?

Myriad of different ways. A useful distinction:

- Gradient–free methods

- Gradient–based methods

# Gradient Free Methods

- Graph–based genetic algorithms

  - Mutations and crossover on a pool of candidates

  - Elitist natural selection rule

- Yoshikawa et. al. propose using SMILES

  - Population of SMILES

  - Grammatical Evolution

- Many more...



Mating Pool     Crossover     Mutation

# Gradient Based Methods

- Recurrent Neural Networks

- (Variational) Autoencoders

- Normalizing Flows

- Generative Adversarial Networks (GANs)

# Recurrent Neural Networks

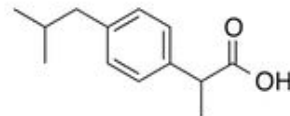- Work on sequences (SMILES)

- Goal: given training sequences → learn to generate new sequences that resemble those of training.



Ibuprofen
CC(C)Cc1ccc(cc1)C(C)C(O)=O

- Sequence: $S_{1:T} = (S_1, \ldots, S_T)$ where $S_i \in \mathcal{V}$

- Training: maximum likelihood, equiv to minimize loss function:

$$L^{MLE} = -\sum_{s \in \mathcal{T}} \sum_{t=2}^{T} \log \pi_\theta(s_t | S_{1:T-1})$$

- Generation: sequentially sample from multinomial dist.
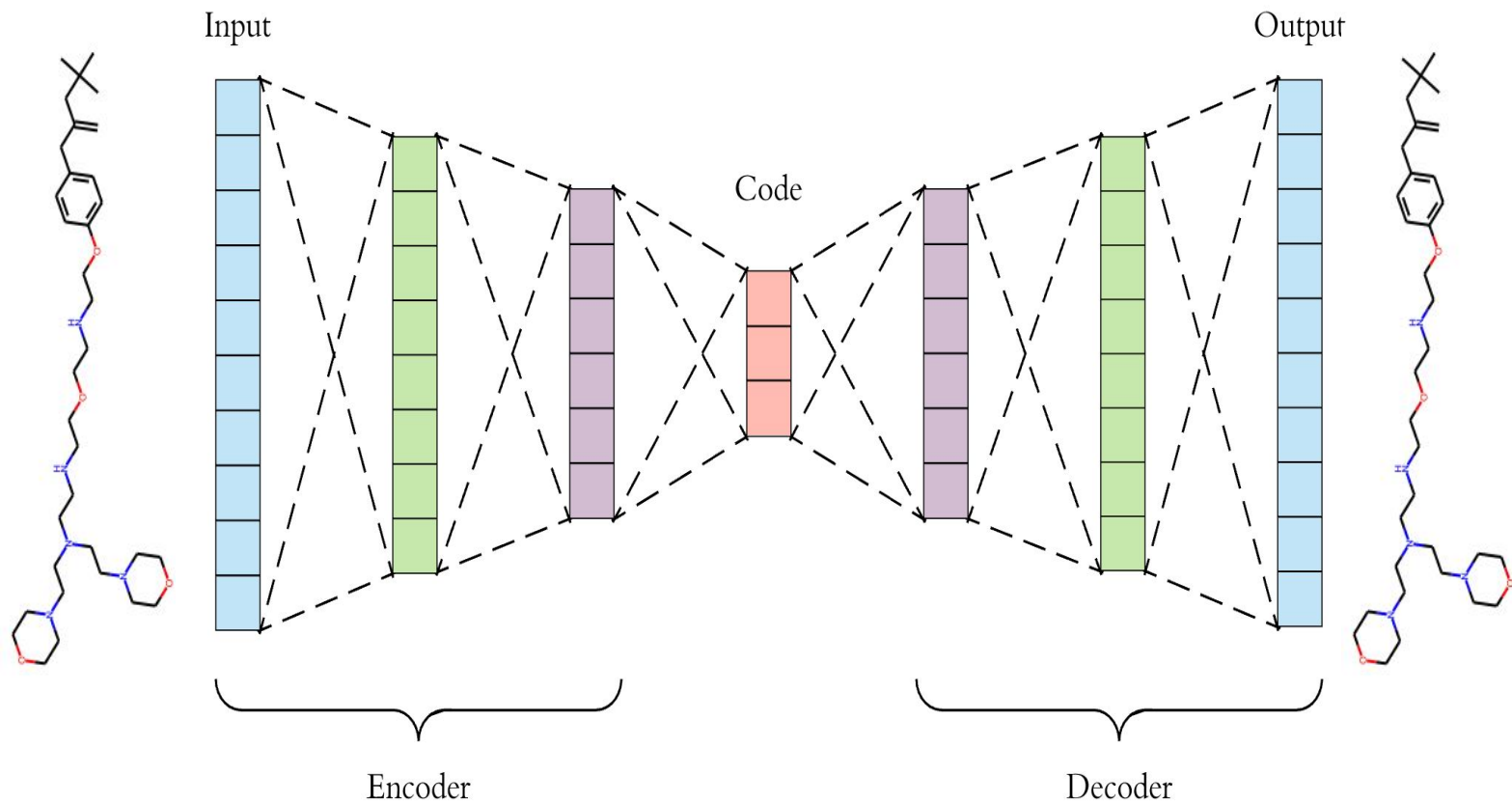
- Thermal rescaling

$$\hat{p}_i \propto \exp\left(\frac{p_i}{T}\right)$$

# (Variational) Autoencoders

# Variational Autoencoders

- Goal: learn probabilistic latent variable model for data generation

$$z \sim p(z)$$
$$x \sim p_\theta(x|z)$$

- We want to maximize $p(x) = \int p_\theta(x|z)p(z)dz$; instead maximize

$$\log p(x) \geq \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right]$$

- RHS is equal to

$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_{\theta(x|z)} \right] - D_{KL}\left[ q_\phi(z|x), p(z) \right]$$

# Variational Autoencoders

- Typically: $p(z)$ independent standard normal dist. and $q_\phi(z|x)$ factorized multivar. normal
- Mean and variance functions of encoder parameterized through CNN.
- Decoder normally RNN

- Training

  - Encode each training sample x into z

  - Decode z into x'

  - Minimize loss function

- Generation

  - Get point in latent space z

  - Decode z sampling $x \sim p_\theta(x|z)$

# Normalizing Flows



- Learn series of parametric bijective transformations of probability distributions

- Allows (easy) calculation of exact likelihood.

- Deep NN with bijective layers

# Generative Adversarial Networks



Generative adversarial network (GAN)

- Generator: generate molecule from Gaussian noise
- Discriminator: distinguish real from fake molecules
- Train to compete against each other

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{x \in p_{\mathrm{d}}(x)} \Big[ \log D(x) \Big]$$
$$+ \mathbb{E}_{z \in p_z(z)} \Big[ \log \big( 1 - D(G(z)) \big) \Big]$$

# Recall that...

*"De novo molecular design is the process of* **automatically proposing novel chemical structures** *that* **optimally satisfy desired properties**"

Combinatorial, black-box, stochastic, multi-objective optimization with black-box constraints

# Generate molecules that optimally satisfy desired properties

- Goal: learn valid molecules with **desirable properties**

- Infeasible to measure properties experimentally for every generated molecule...

- Infeasible to use computational chemistry to compute properties...

- **Prediction**: quantitative structure–activity relationship (QSAR)

- Done usually in separate datasets

- Many models depending on property, representation, etc.

    - Molecular Descriptors
    - SMILEs
    - Graphs

# Using properties to guide generation

1.  Reinforcement Learning coupled with sequence generator

    - A time t, state is $(s_0, \ldots, s_t)$
    - Action is next token $a_t = s_{t+1}$
    - After taking action, a reward $R_t$ is perceived
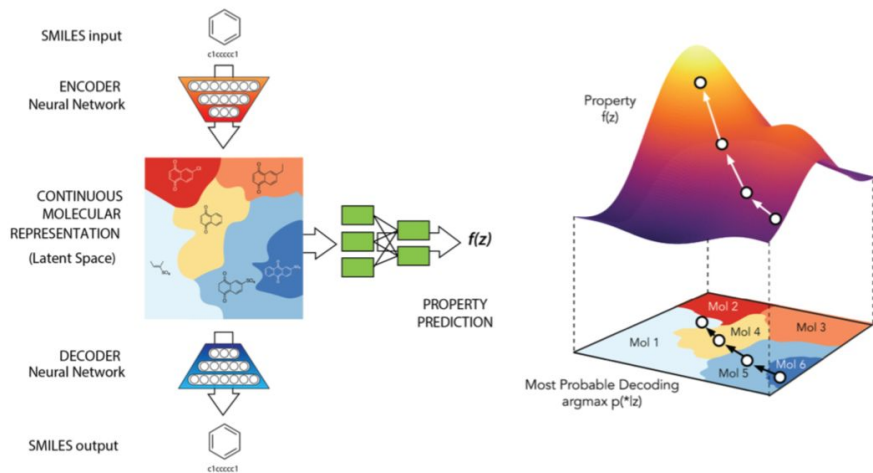    - Goal, learn policy $\pi_\theta(a|s)$

    $$\max_\theta \mathbb{E}[\sum_{i=1}^{T} R_i | s_0, \theta]$$

    - The only non-zero reward is $R_T$ which is equal to the property prediction

# Using properties to guide generation

2. Optimization with VAE

- Learn map from latent space to property (e.g. through GP)
- Optimize that map (gradient ascent, bayesian optimization, etc.)

# Issues/Thoughts

- Multi-objective optimization

  - Many properties to be optimized (depending even on different stakeholders!)

  - Drug discovery: **high binding affinity to biological target**, low toxicity, solubility, synthetically accessible, stability, economical costs!

  - Commonly: predict properties independently and combine predictions in loss function.

  - Also, hold properties constant implicitly through structural constraints.

  - **Decision theory**: **multi-attribute utilities** to incorporate different objectives for different stakeholders into the generative process

# Issues/Thoughts

- Uncertainty quantification

    - Models rely on predictions to generate promising molecules

    - Accuracy of these models is key

    - In small data regimes... models tend to be less accurate.

    - Incorporate uncertainty quantification into generative process! (**Bayesian inference**)

    - Exploration vs exploitation (**Bayesian optimization**)

    - **Bayesian decision theory**
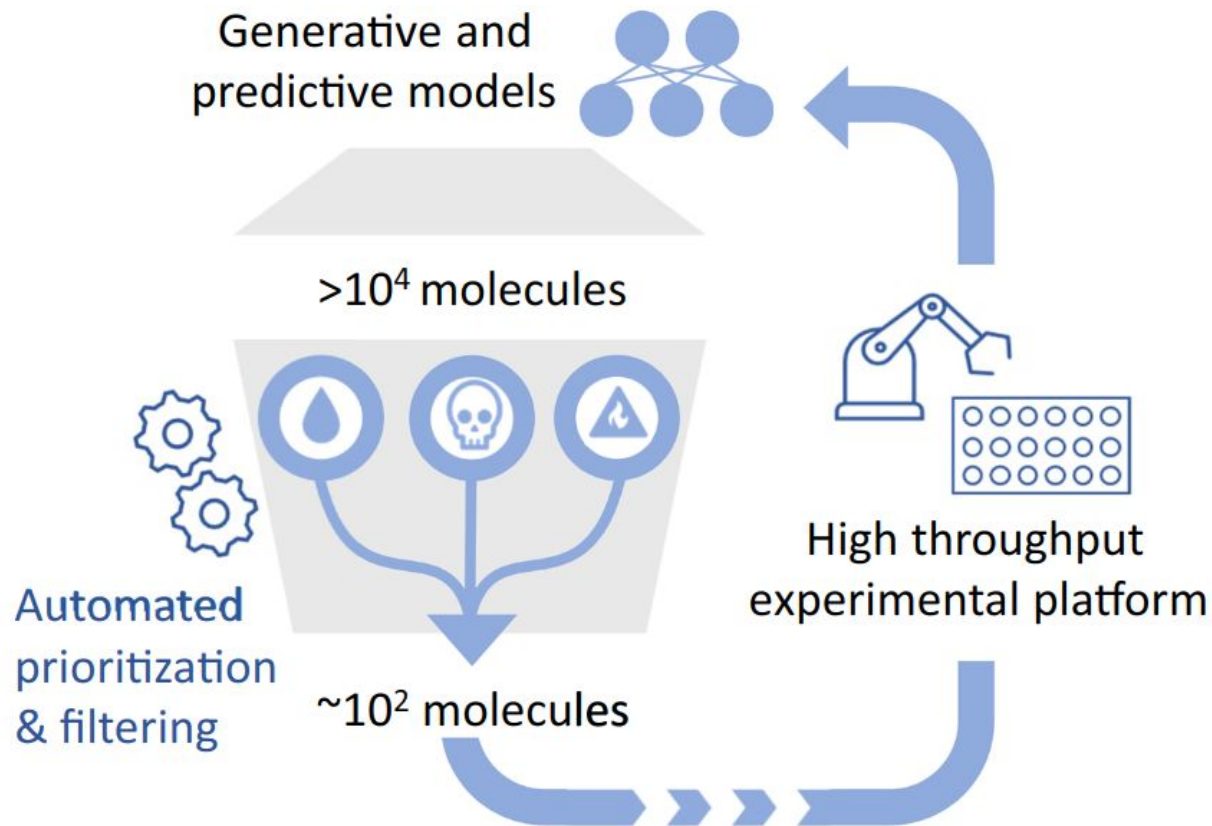
# Issues/Thoughts

- Synthesizability

  - Generated molecules must be easy to synthesize

  - This concept is hard to define!

  - Methods to automatically evaluate synthesizability without human intervention

  - Rather than molecules, generate synthetic pathways  (learn reactions)

# Other relevant fields

- **Graph based deep learning**

- **Geometric deep learning**

- **Combinatorial black-box optimization**

- **Heuristic search algorithms**

- **Reinforcement Learning**

Generative and predictive models

$>10^4$ molecules

Automated prioritization & filtering

$\sim10^2$ molecules

High throughput experimental platform

# The reality?

- More likely: computer–aided molecular design

- Interpretability

    - Prediction is not enough, we need understanding (?).

    - Chemist need to derive an actionable hypothesis from model output.

    - If chemist sees, e.g. structural elements responsible for toxicity, she might have ideas on how to modify molecule to diminish toxicity

    - Interpretable representations: molecular descriptors…?

    - Interpretable methods to determine **causality between structure presence and property** (**causal inference, counterfactual inference**)

# References

Bilodeau, C., Jin, W., Jaakkola, T., Barzilay, R., & Jensen, K. F. (2022). Generative models for molecular discovery: Recent advances and challenges. Wiley Interdisciplinary Reviews: Computational Molecular Science, e1608.

Meyers, J., Fabian, B., & Brown, N. (2021). De novo molecular design and generative models. Drug Discovery Today, 26(11), 2707–2715.

Elton, D. C., Boukouvalas, Z., Fuge, M. D., & Chung, P. W. (2019). Deep learning for molecular design—a review of the state of the art. Molecular Systems Design & Engineering, 4(4), 828–849.

Gallego, V., Naveiro, R., Roca, C., Ríos Insua, D., & Campillo, N. E. (2021). AI in drug development: a multidisciplinary perspective. Molecular Diversity, 25(3), 1461–1479.

Yoshikawa, N., Terayama, K., Sumita, M., Homma, T., Oono, K., & Tsuda, K. (2018). Population-based de novo molecule generation, using grammatical evolution. Chemistry Letters, 47(11), 1431–1434.
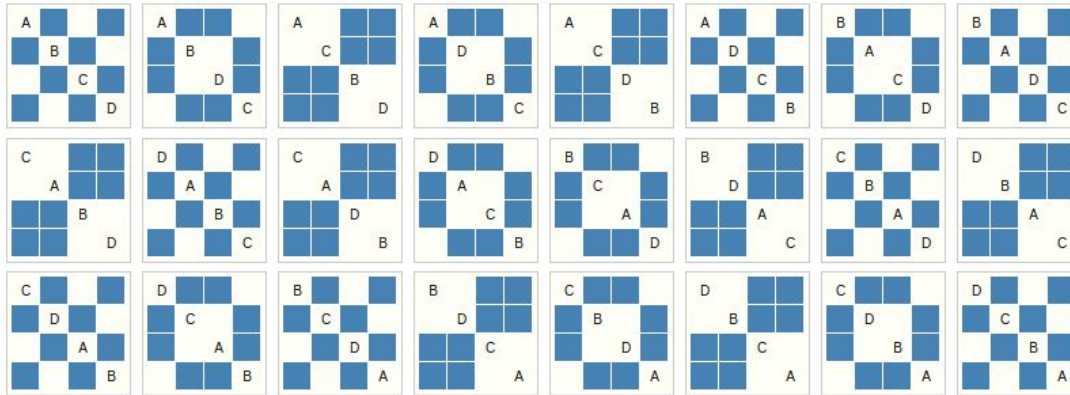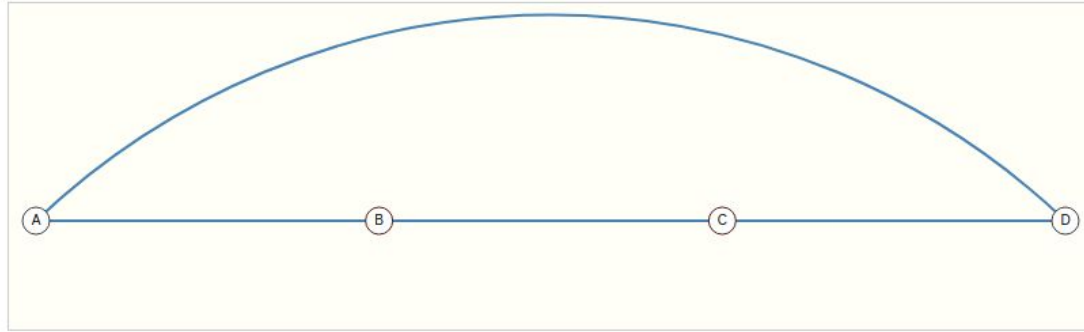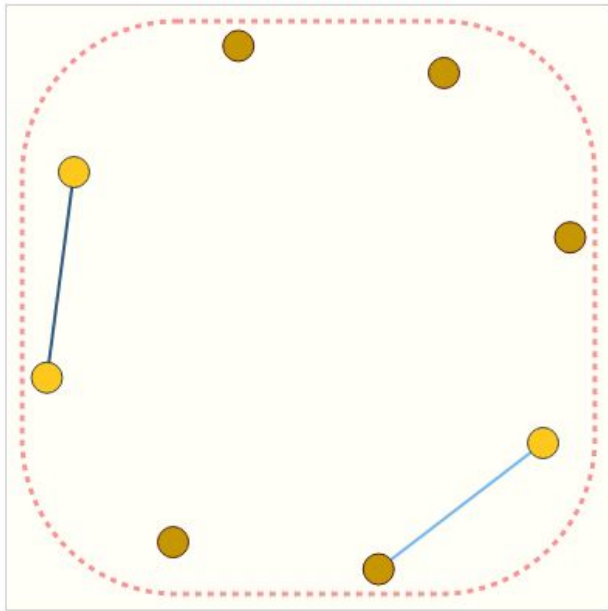
# Thanks!



roi.naveiro@icmat.es

https://roinaveiro.github.io/

https://github.com/roinaveiro

# Adjacency Matrices

# Permutation Invariant representation

# Unconstrained generation

- Goal: learn general distribution of molecules in chemical space

- Evaluated based on chemical validity, novelty, uniqueness