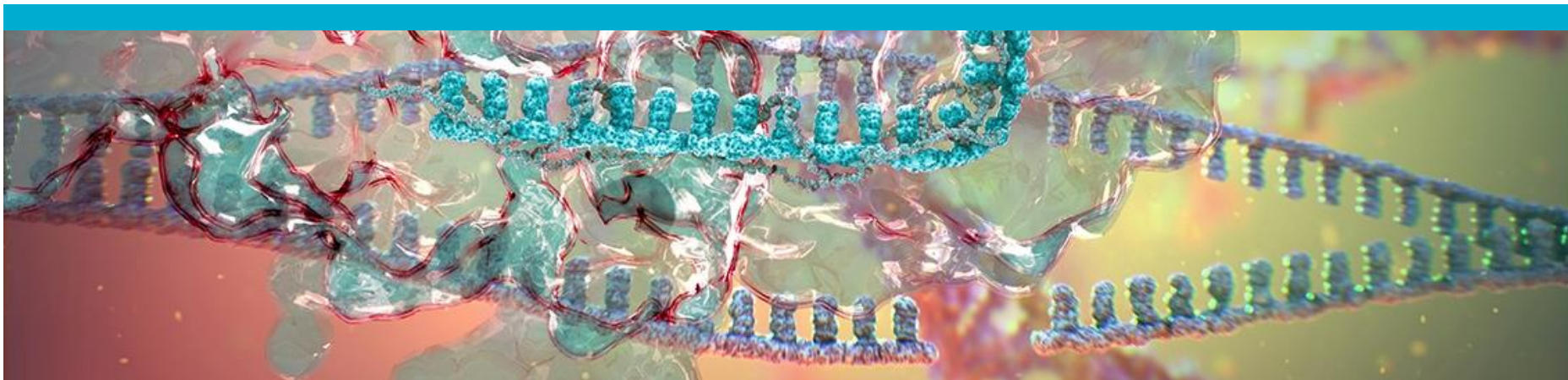# AI for drug design an industrial perspective

**Ola Engkvist, Molecular AI, Discovery Sciences, R&D, Gothenburg, Sweden**
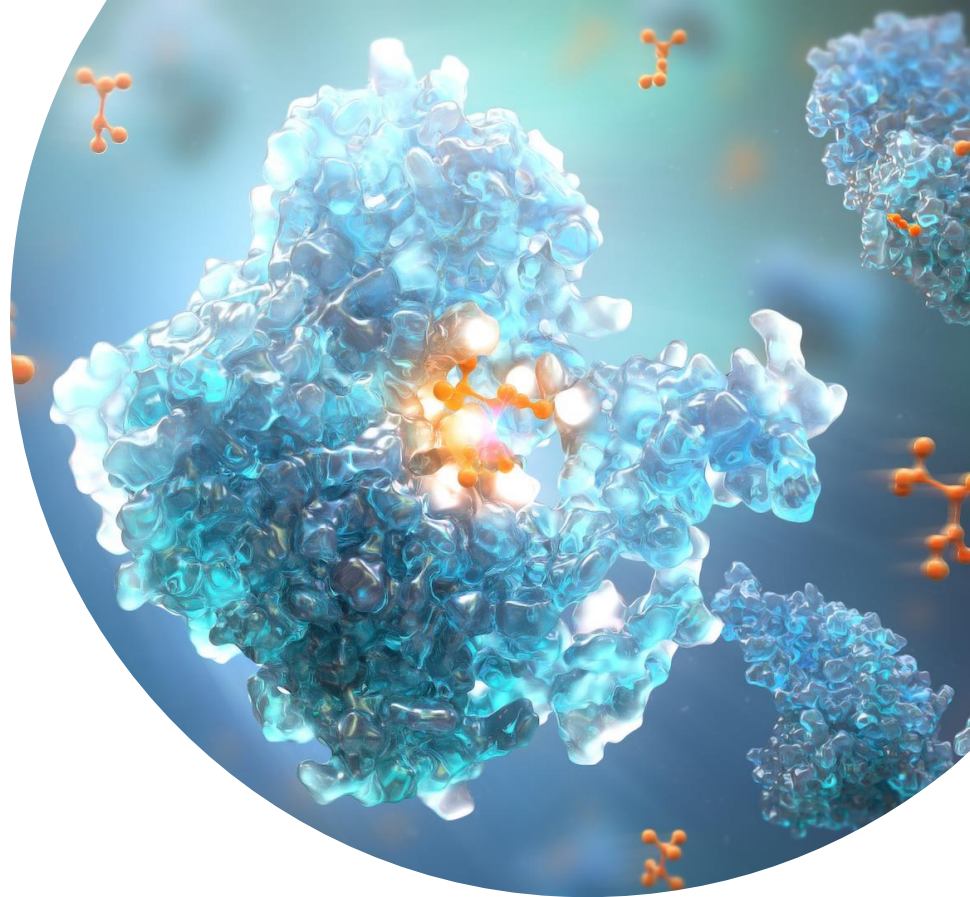
Chemoinformatics Strasbourg Summer School 2022
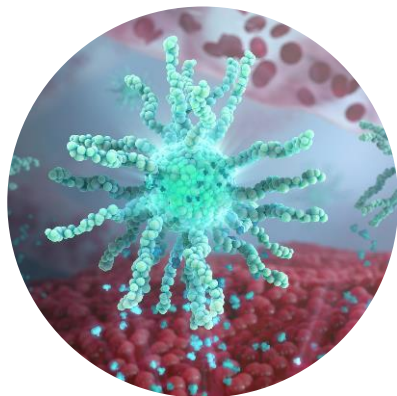
July 1 2022

# We push the boundaries of science to deliver life-changing medicines

Inspired by our **values** and **what science can do**, we are **focused** on **accelerating the delivery of life-changing medicines** that create enduring **value for patients and society.**

# 2021 global dimensions

**$37.4bn**
Total Revenue (incl. COVID-19 vaccine)

**+38%**
Total Revenue growth (23% excl. COVID-19 vaccine)

**$9.7bn**
invested in our science

**13**
medicines with annual sales of more than $1 billion

**$13bn**
Oncology Product Sales

**$8bn**
Cardiovascular, Renal & Metabolism Product Sales

**$6bn**
Respiratory & Immunology Product Sales

**$3bn**
Rare Disease Product Sales (from 21/7/2021)

**22**
Regulatory approvals and authorisations in major markets

**2.5bn**
COVID-19 vaccine doses supplied to more than 180 countries together with our partners

**110**
successful markets launches

**161**
projects in clinical phase of development

**83,100**
employees (Dec. 2021)

**87%**
of employees believe strongly in our future direction and key priorities (Nov. 2021)

**59%**
Reduction in Scope 1 and 2 greenhouse gas emissions since 2015

**31m**
people reached through our Access to Healthcare programmes

NOTE: All growth rates at Constant Exchange Rates
Source: 2021 Annual Report

# Global, science-led, patient-focused biopharmaceutical company

Science and innovation-led

Therapy areas of focus: Oncology; Cardiovascular, Renal & Metabolism; Respiratory & Inflammation; Rare Disease

Diversified portfolio with broad coverage across primary care, specialty care and rare diseases

Commitment to people and society

Global strength, with balanced presence across regions

# Our Strategic Priorities

Deliver growth and therapy area leadership

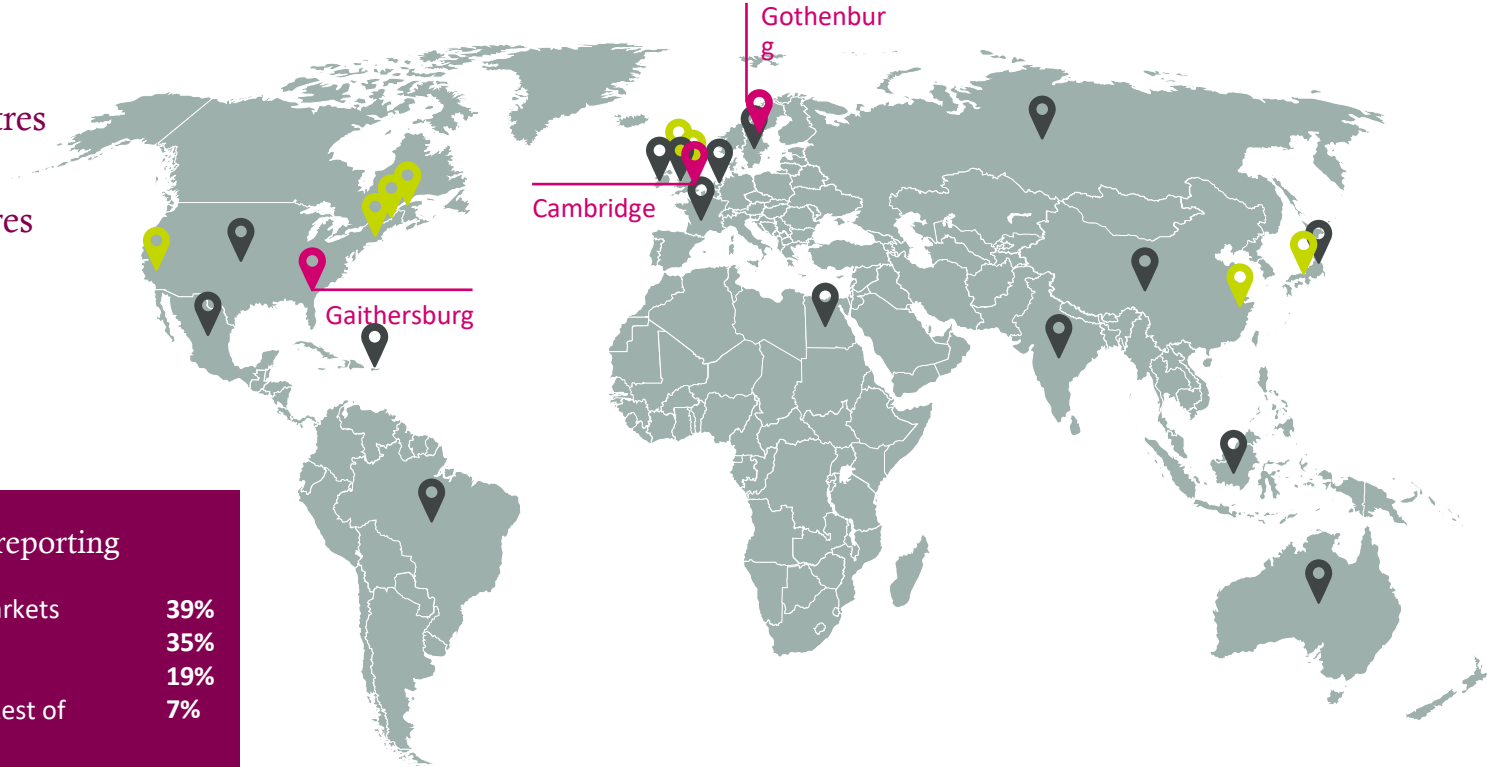Accelerate innovative science

Be a great place to work

# Global reach and presence

📍 3 global R&D centres

📍 8 other R&D centres and offices

📍 28 manufacturing sites in 16 countries

Gothenburg

Cambridge

Gaithersburg

## Employees by reporting region

| | |
|---|---|
| Emerging Markets | 39% |
| Europe | 35% |
| US | 19% |
| Established Rest of World | 7% |

# AstraZeneca Gothenburg

# AstraZeneca in Gothenburg

Gothenburg has a unique culture of **collaboration and open innovation**, supporting all our main therapy areas throughout the entire life cycle of our medicines.



Gothenburg



**Supporting the entire life cycle of medicines:** drug discovery, advanced drug delivery, medical device development, manufacturing for clinical trials.

**~2,800 employees** of **70+ nationalities,** including 600 PhD researchers and 30 professors. 57.9 percent of our line managers in Gothenburg are women.

**Reduced the site's carbon dioxide emissions by 99% since 2015** with a clear vision to further reduce our carbon footprint.
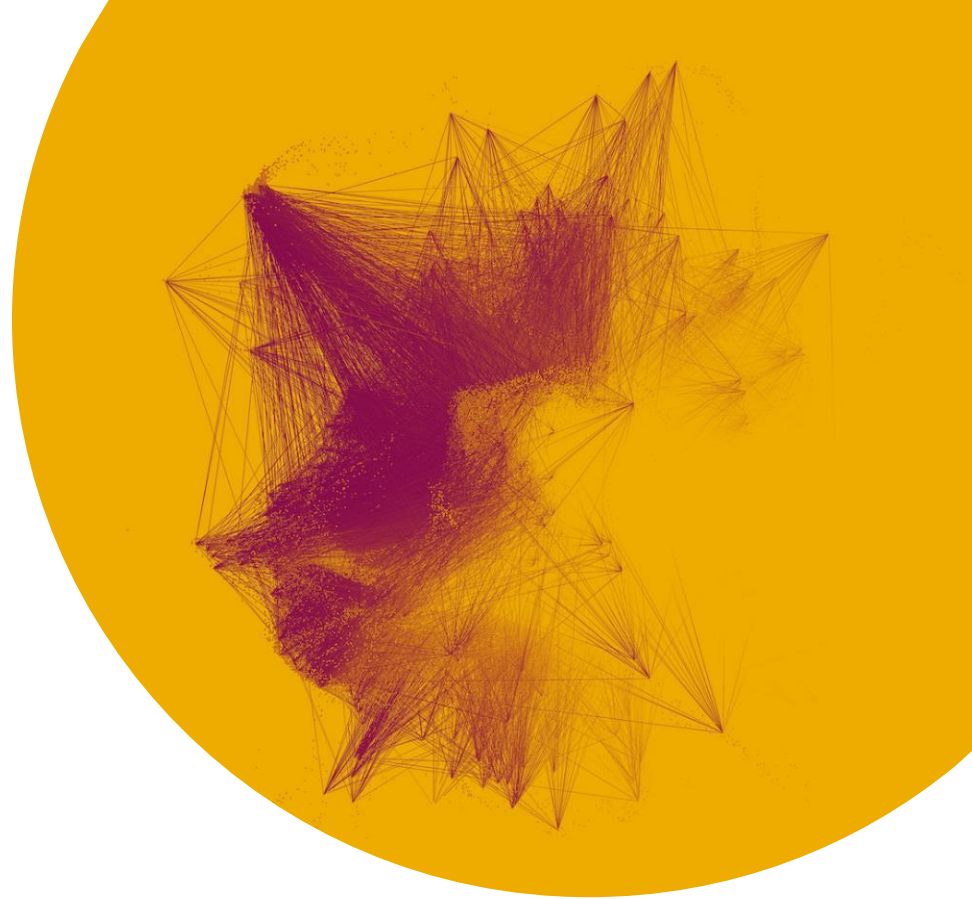
# From idea to patient

One of the site's unique features is that we have most of the resources and expertise needed to support the entire lifecycle of a medicine. That's everything from idea generation, via clinical development, to pilot scale manufacturing and distribution, and global commercialisation and product maintenance.



Inputs
> Applying our resources to meet unmet medical need

Outputs
> Returns to shareholders
> Improved health

Research and development phases 10–15 years

1 Find potential medicine

2 Pre-clinical studies

3 Phase I studies

4 Phase II studies

5 Phase III studies

6 Regulatory submission and pricing

7 Launch new medicine

8 Post-launch research and development

9 Post-exclusivity

Post-exclusivity 20+ years

Launch phase 5–10 years

Our Purpose

# Data science & AI: transforming drug discovery and development

AI and machine learning are transforming the way we discover and develop new medicines.

We aim to increase the probability of success and reduce drug discovery and development timelines by applying advanced AI and machine learning across R&D.

# Where can AI impact drug discovery and development

# Drug Design

**Which compound to make next?**

**How to make the compound?**

# The Design Make Test Analyze cycle in Drug Design

**Drug target**

**Chemical starting point ("Hit") found through HTS, DEL, fragment screening or knowledge**

- Weakly active
- Target unselective
- Toxicity risk
- Low metabolic stability

**Design**

**Make**

**Analyse**

**Test**

**Candidate drug**

- Highly potent
- Effective in *in vivo* models
- Metabolically stable
- No toxicity issues

**~3 years**

**Multiple of DMTA cycles**

# AI based drug design
## How can we reduce the time to deliver a clinical candidate?



Select the most efficient synthetic route

Design

Analyze

Make

Test

Make information rich compounds in each cycle

Increase speed

Maximize learning

# Why now?

Why would this presentation have been science fiction 5 years ago?

➢ Increased computational power

    Never underestimate an exponential law

➢ Advances in neural network algorithms

    New algorithms in other fields that can be adapted to our needs i.e. Image recognition, <u>Natural language processing</u>, Playing Go

➢ Open-source software

    Python, RDKit, scikit-learn, PyTorch, Tensorflow

How can we take advantage of the progress in Natural Language Processing?

Molecules can be described with the language SMILES



Language Translation ⟷ Synthesis prediction

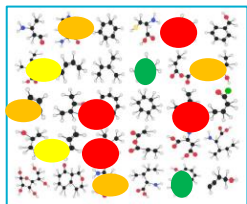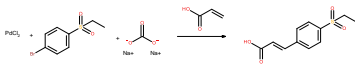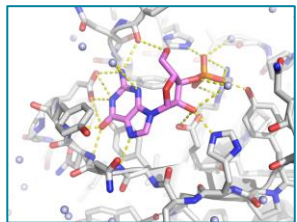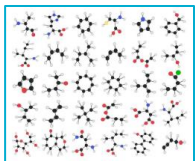Language Translation ⟷ Molecular optimization

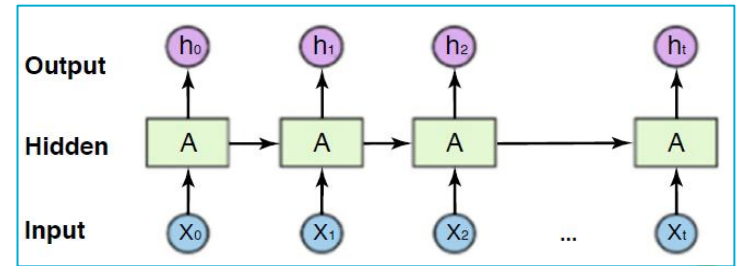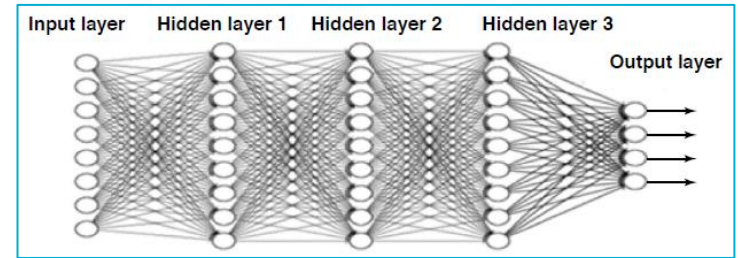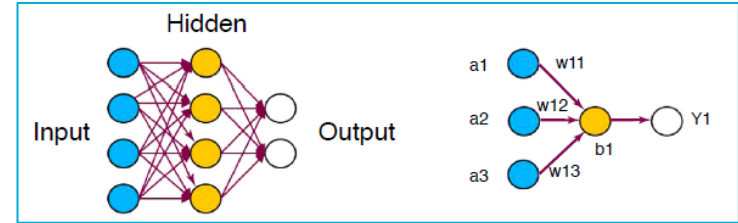Text generation ⟷ Chemical space exploration

# Where are we with AI based drug design?



- ✓ Deep learning based molecular de novo generation
  - ✓ Algorithms to navigate the whole relevant chemical space now exists
  - ✓ Scoring the generated molecules the main bottleneck

- ✓ Synthetic route prediction
  - ✓ Powerful new algorithms have been developed
  - ✓ Further progress needs better data

- ✓ Molecular property prediction
  - ✓ Novel more flexible deep learning based methods
  - ✓ No AF2 moment, No progress in prospective competitions (SAMPL, IDG, Cache)

- ✓ 3D protein prediction
  - ✓ Stunning progress with AlphaFold2
  - ✓ Dynamics needs to be included for major impact

- ✓ It is not only about ML/AI based technology
  - ✓ AI+ vision need to include high-throughput data generation, automation & physics-based modelling
  - ✓ AI first culture
  - ✓ Continiuous training and education of staff

# Neural Networks & Deep Learning

- **Neural Networks known for decades**
  - Inputs, Hidden Layers, Outputs
  - Single layer NNs have been used in QSAR modelling for years
- **Recent Applications use more complex networks such as**
  - Multi-layer Feed-Forward NNs
  - Convolutional NNs
    - biological image processing
  - Auto-encoder NNs
  - Recurrent NNs
    - Trained using Maximum Likelihood Estimation to maximize the likelihood of next character
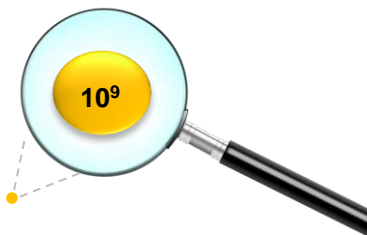
# Generative Model vs Enumeration for molecular discovery

## Physical Storage Size

## Size of Molecular Space

**Traditional Enumeration**

41 GB

$10^9$

**Generative Model**

50 MB

$10^{60}$

**Generative models can sample practically unlimited chemical space**

**Generative models do not contain any explicit molecules but generate them probabilistically**

# Two different ways how can AI help finding the next molecule to make?



Hit Finding & scaffold hopping
Sample the whole chemical space

## Recurrent Neural Networks



Molecular Optimisation
Sample a focused chemical space



## Transformer

# Recurrent Neural Network & Natural language generation

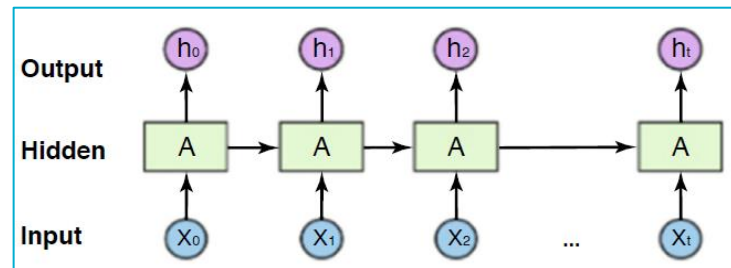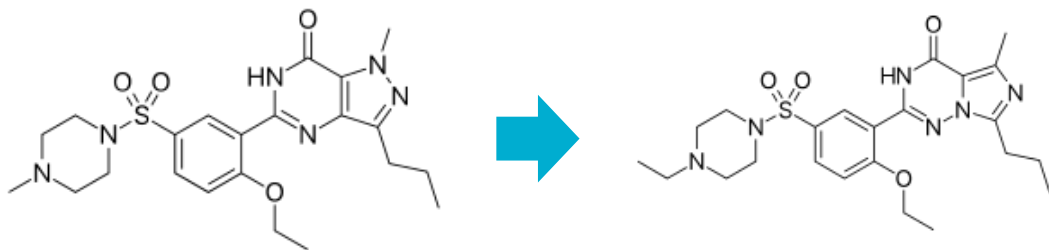# Natural language generation and molecular structure generation

- Can we borrow concepts from natural language processing and apply to SMILES description of molecular structures to generate molecules?

$$\text{The} \longrightarrow \text{grass} \longrightarrow \text{is} \longrightarrow ?$$

- Conditional probability distributions given context

- $P(green \mid is, grass, The)$

$$C \longrightarrow C \longrightarrow = \longrightarrow ?$$

- $P(O \mid =, C, C)$

# Tokenization of SMILES

- Tokenize combinations of characters like "Cl" or "[nH]"

- Represent the characters as one-hot vectors

Graph:

SMILES:                    ClCc1c[nH]cn1

One-hot encoding:

|       | Cl | C | c | 1 | c | nH | c | n | 1 |
|-------|----|---|---|---|---|----|---|---|---|
| **C** | 0  | 1 | 0 | 0 | 0 | 0  | 0 | 0 | 0 |
| **c** | 0  | 0 | 1 | 0 | 1 | 0  | 1 | 0 | 0 |
| **n** | 0  | 0 | 0 | 0 | 0 | 0  | 0 | 1 | 0 |
| **1** | 0  | 0 | 0 | 1 | 0 | 0  | 0 | 0 | 1 |
| **nH**| 0  | 0 | 0 | 0 | 0 | 1  | 0 | 0 | 0 |
| **Cl**| 1  | 0 | 0 | 0 | 0 | 0  | 0 | 0 | 0 |

# The generative process



Sampled SMILES · Log P · Structure

# AI live: Create Structures Similar to Celecoxib



- **Key Message**
  - RNN generates structures similar to Celecoxib
  - Rapid sampling!
  - Average score describes how many learning steps are required to reach similar compounds

# To think about when using reinforcement learning

- RL will exploit loopholes in the scoring function
- RL will exploit the first minima it finds



Scaffold penalty to assure diverse scaffolds are identified

**Blaschke et al Journal of Cheminformatics 2020**

# Docking is essential to score molecules



**DockStream: A Docking Wrapper to Enhance De Novo Molecular Design**

**Guo et al ChemRxiv**

# Our knowledgebase: ReactionConnect



Predictive Reaction Models

**iLab, MedChem, PharmDev**

# AiZynthFinder

Web-GUI based on MIT MLDPS consortium tools

Scripting access via Python Objects

Jupyter based GUI

# So which lessons have we learned?

- Needs in discovery chemistry and process chemistry are very different

- Extracting and integrating reaction data is hard work

- Challenging to assess the utility of different tools

- Advanced building block look-up

- Impact on synthetic routes have mainly been from specialized tools like the Ringbreaker

# Artificial Intelligence Guided Drug Design Platform (REINVENT)

## Molecular Generation



- Recurrent neural networks
- Reinforcement learning
- Scaffold penalty
- Library design

## Scoring



- Synthetic accessibility score
- QSAR models (ADME, off-target)
- 3D shape
- Docking

## Postprocessing



- free energy perturbation binding affinity
- Synthetic route assesment
- Clustering

**Core is Based on Open Source Software**

**Commercial Plugins when appropriate for scoring**

# AI+ vision for drug design

## AI can't transform drug design alone

### High-Throughput Data Generation



- The most important determinant of the usefulness of a model is the size and quality of the data set for training

- High-Throughput Experimentation for generating chemical reaction data

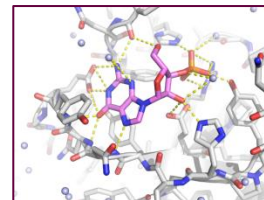- Cell-paint & transcriptomics to create molecular signatures

- DNA Encoded Library models to score molecules

### Automatize Make & Test



- Autonomous optimization of compounds is needed to radically cut timelines for clinical candidate delivery

- Multistep reactions with intermediate purification on automation platform

- Automatic testing after synthesis & purification

- Autonomous decision making under uncertainty which compounds to make

- Human-in-the-loop modelling

### Combine AI with physics



- More accurate models for difficult to predict properties can be created through combining physics and AI

- Relative binding free energy perturbation binding affinity in molecular optimization

- Absolute binding free energy perturbation to estimate binding energies in hit finding and for scaffold hopping

- Estimation of thermodynamic solubility
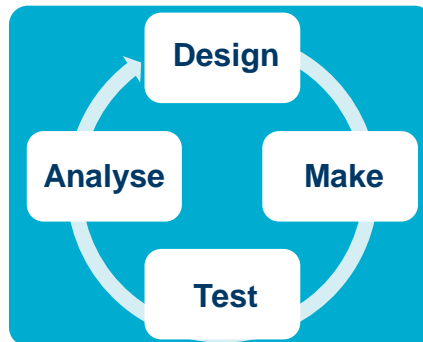
- Combine ML/MD to identify cryptic pockets

# Integration of AI and automation



**AI**

**iLAB**

# Automated Synthesis Platform @AstraZeneca

# Keep a balanced view!

➢ **Progress have been made!**

  ➢ Molecular Generation

  ➢ Synthesis prediction

  ➢ AlhpaFold2

➢ **Progress will continue!**

  ➢ Better hardware

  ➢ High-throughput data generation

  ➢ Novel innovative algorithms

  ➢ Better deep learning based force-fields for affinity estimation etc

➢ **Some predictions will be difficult!**

  ➢ In vivo properties from only molecular structures

  ➢ Expensive data generation, Noisy data, Non-smooth surface to learn

# Why am I an optimist

- Increased computational power
  - Never underestimate an exponential law

- Increased automation provides large and consistent datasets
  - HTE, DEL, Cell Paint, RNA-seq

- Advances in computational algorithms
  - Merging of physics-based modelling & ML
  - GPT-3, Codex

# What about AlphaFold2?

➢ Terrific achievement!

  ➢ Winning a prospective competition with margin based on public data!

  ➢ Big Science (People, Compute)

  ➢ Public release will encourage further development & innovation

  ➢ Looking forward to the next generation of models (capturing protein dynamics, RNA structures)

➢ Impact on drug design

  ➢ Facilitate solving x-ray and Cryo-EM structures

  ➢ Lack of protein dynamics have limited the use so far

# What does success look like?

➤ Metrics like time saving are the results of success not the success itself

➤ Trust in the AI designed molecules in the same way as for instance x-ray crystal structures are trusted

- ➤ Trust in the predictions for individual molecules
- ➤ Trust that the AI generated molecules are the best molecules taking the project most efficiently to a clinical candidate

# What are the challenges for AI driven drug design?

- Scaling ML/AI solutions for drug design to a whole drug discovery project portfolio including projects with low data volume
    - (pre-trained) molecular transformers
    - Privacy-preserving ML

- Physics based modelling
    - Binding affinity and solubility predictions are major bottlenecks

- "Cambrian revolution" of new AI methods makes it difficult to assess progress

- Flexibility of chemistry automation

- Educational, cultural & logistical challenges besides scientific

# Science Molecular AI @AZ



**ACS Central Science** — Research Article
Cite This: ACS Cent. Sci. 2018, 4, 120–131

**Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks**

**RESEARCH**

Molecular De-Novo Design through Deep Reinforcement Learning

Marcus Olivecrona*, Thomas Blaschke[†], Ola Engkvist[†] and Hongming Chen[†]

**RESEARCH ARTICLE** — Open Access

Exploring the GDB-13 chemical space using deep generative models

Josep Arús-Pous[1,3]*, Thomas Blaschke[1,4], Silas Ulander[2], Jean-Louis Reymond[3], Hongming Chen[1] and Ola Engkvist[1]

**JCIM** — JOURNAL OF CHEMICAL INFORMATION AND MODELING

pubs.acs.org/jcim — Application

**REINVENT 2.0: An AI Tool for De Novo Drug Design**

Thomas Blaschke, Josep Arús-Pous, Hongming Chen, Christian Margreitter, Christian Tyrchan, Ola Engkvist, Kostas Papadopoulos, and Atanas Patronov*

**Journal of Medicinal Chemistry**

This is an open access article published under an ACS AuthorChoice License, which permits copying and redistribution of the article or any adaptations for non-commercial purposes.

pubs.acs.org/jmc — Article

**"Ring Breaker": Neural Network Driven Synthesis Prediction of the Ring System Chemical Space**

Amol Thakkar,* Nidhal Selmi, Jean-Louis Reymond, Ola Engkvist, and Esben Jannik Bjerrum*

**Chemical Science** — ROYAL SOCIETY OF CHEMISTRY

**EDGE ARTICLE** — View Article Online

Cite this: Chem. Sci., 2021, 12, 3339

All publication charges for this article have been paid for by the Royal Society of Chemistry

**Retrosynthetic accessibility score (RAscore) — rapid machine learned synthesizability classification from AI driven retrosynthetic planning†**

Amol Thakkar,[*ab] Veronika Chadimová,[a] Esben Jannik Bjerrum,[a] Ola Engkvist[a] and Jean-Louis Reymond[*b]
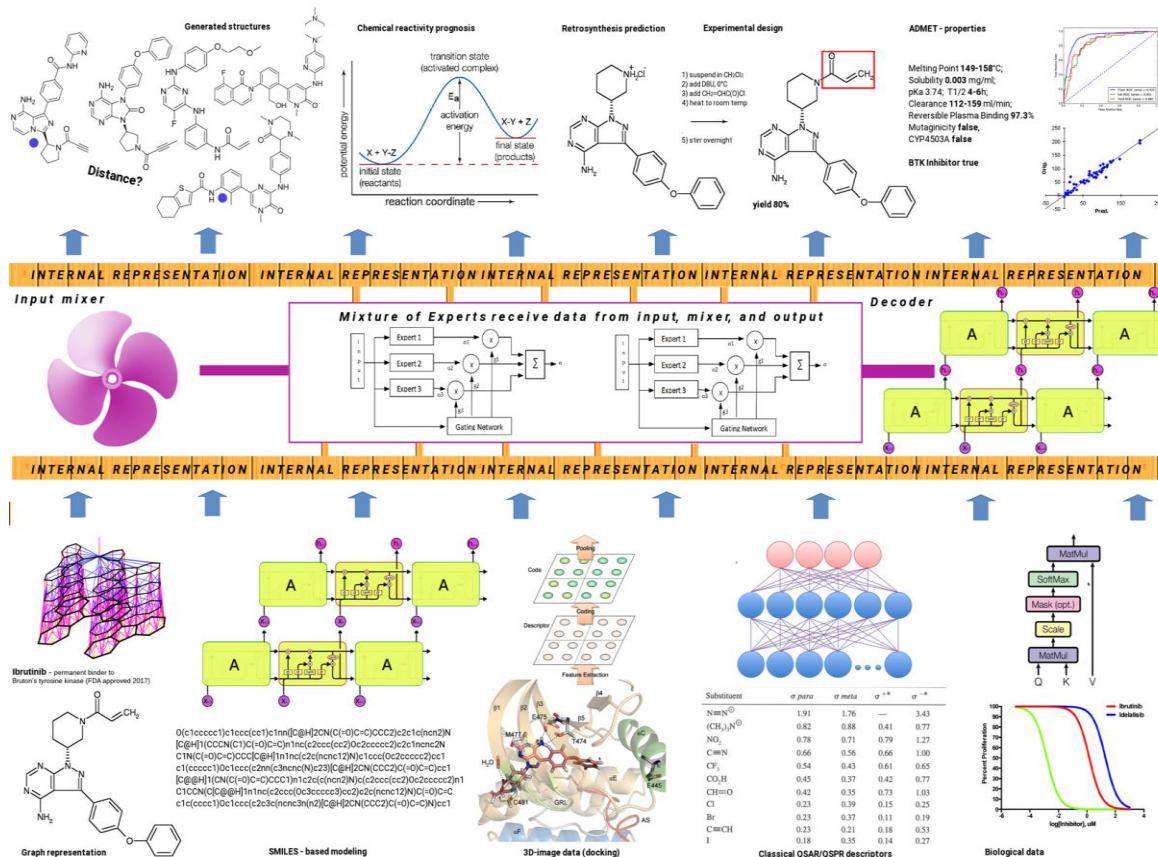
**SOFTWARE** — Open Access

**AiZynthFinder: a fast, robust and flexible open-source software for retrosynthetic planning**

Samuel Genheden[1*], Amol Thakkar[1,2], Veronika Chadimová[1], Jean-Louis Reymond[2], Ola Engkvist[1] and Esben Bjerrum[1*]

**40**

Open Source: https://github.com/MolecularAI

# Advanced Machine Learning for Innovative Drug Discovery https://ai-dd.eu



✓ On-line courses
✓ Schools in AI & ML

Public lectures by ZOOM

See announcements at:

https://ai-dd.eu/news

Or follow:

https://twitter.com/**AiddOne**