# Natural Language Processing for Drug Discovery: The State of Practices, Opportunities, and Challenges

**Suneel Kumar BVS, Ph.D.,**

Director of Drug Design & AI
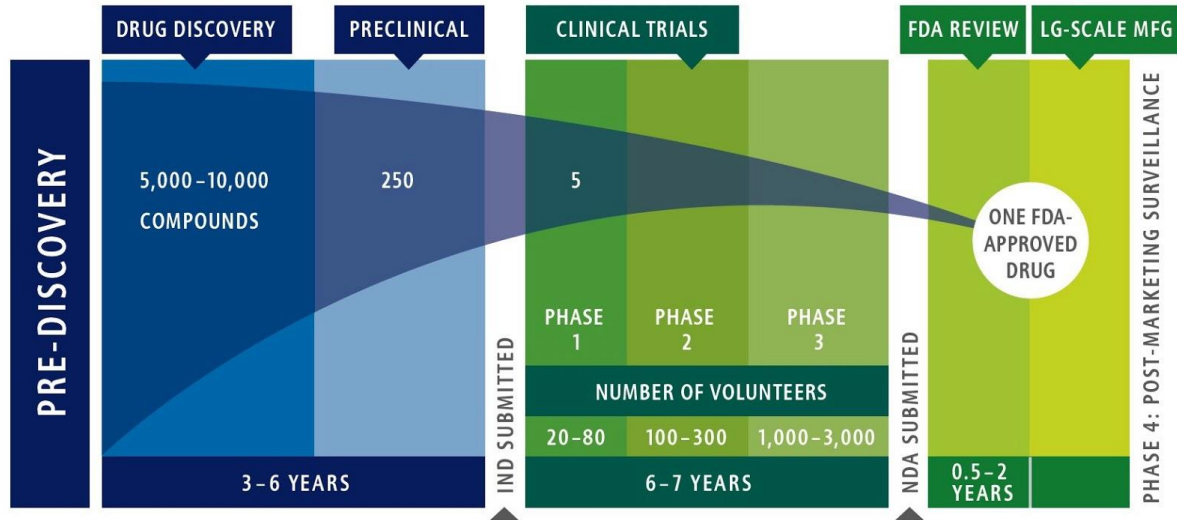Molecular Forecaster

# NLP for DD: The State of Practices, Opportunities, and Challenges

Hit identification is a crucial step in the drug discovery process, where potential drug candidates are identified from a large chemical space. There are several methods for hit identification, including structure and ligand based-virtual screening, fragment-based drug design, and AI-driven drug design approaches.

In the context of de novo drug design, both generative methods and NLP are used to extract the information about existing data points, generate and optimize molecular structures based on desired properties. These algorithms can be trained on a variety of data sources, including molecular data, its biological profile, leading to a more informed hit generation process. In this talk, I will share an overview of these methods, practices, limitations, and case studies.

# Drug Discovery - Workflow



Average time/cost for designing one drug = 10 years + $2.6B

# Drug Discovery - Workflow

**Traditional drug R&D takes >10 years and >$2B***

From the discovery to the launch of a new drug



**Hit to Clinical Candidate:**

Cost Involved     : 824$ Millions

Time                : 5.5 Years

Failure Ratio       : 69-85%
(range)

| | Target-to-hit | Hit-to-lead | Lead optimization | Preclinical | Phase I | Phase II | Phase III | Submission to launch | |
|---|---|---|---|---|---|---|---|---|---|
| 1-5 % | 80% | 75% | 85% | 69% | 54% | 34% | 70% | 91% | p (TS) |
| 1-10 | 1.0 | 1.5 | 2.0 | 1.0 | 1.5 | 2.5 | 2.5 | 1.5 | Cycle time (years) |
| $0-1000 | $94 | $166 | $414 | $150 | $273 | $319 | $314 | $48 | Cost per launch (capitalized) $ Millions, 2010 data |

$1.8B Development cost in 2010 (after target discovery)

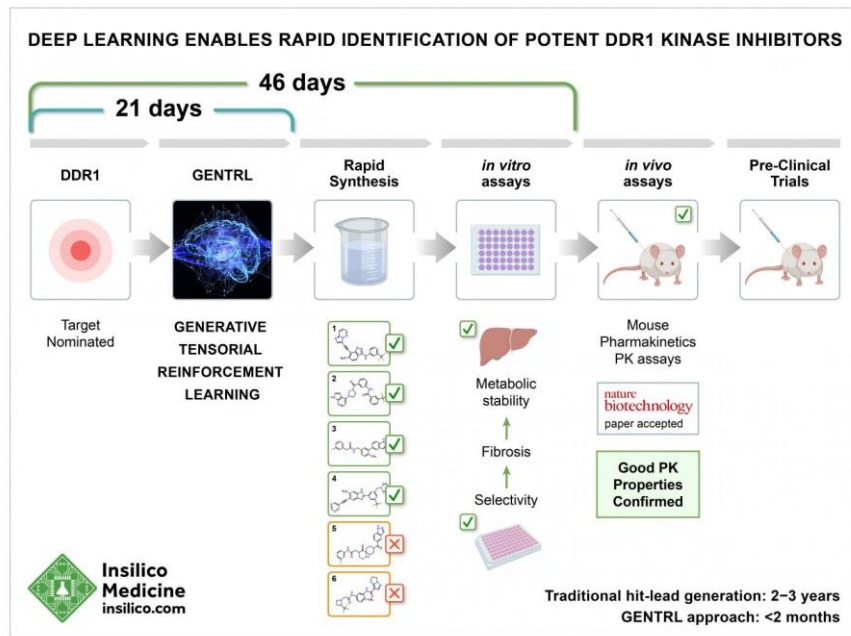**Discovery**        **Development**        **Sales and marketing**

* Modified by Alex Zhavoronkov, PhD, Insilico Medicine from Paul et al, How to improve R&D productivity: the pharmaceutical industry's grand challenge.
Nature Reviews Drug Discovery , 2010
** Based on interviews with the pharmaceutical industry executives

**NLP SUMMIT**

# Insilico Medicine –

## AI generated Lead Candidate



Figure source: Insilico Medicine

# Insilico Medicine –
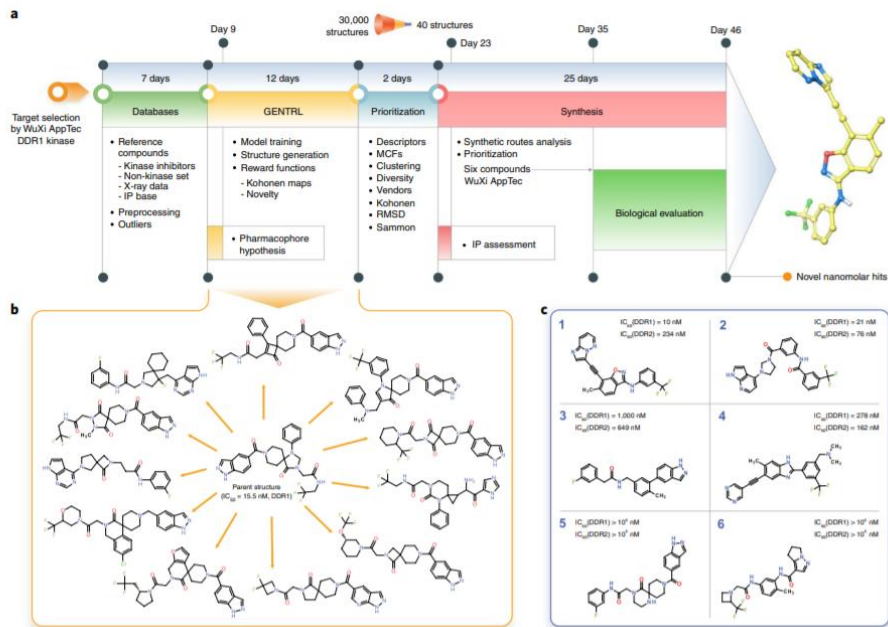
## AI generated Lead Candidate



Fig. GENTRL model design, workflow, and nanomolar hits. a, The general workflow and timeline for the design of lead candidates against human DDR1 kinase.

# Exscientia

## – AI Clinical Candidate

### 1. EXS21546 – Selective, Low-CNS Penetrant A2a Antagonist

| Property/assay | Data |
|---|---|
| **CNS Penetration** | |
| Brain penetration / tumour penetration | Low (Kp, uu brain 0.046 tumour 1.6) |
| **Target binding affinity** | |
| SPR human A2A KD (nM) | 4 |
| Mouse/rat/dog/cyno SPR A2A KD (nM) | 7 / 10 / 63 /3 |
| **Avoid off-targets** | |
| A1/A2B/A3 binding ratio to A2A | 875 / 577 / >1000 |
| **Cell potency** | |
| HEK-Human A2A IC50 (nM) | 37 |
| Human A2A T cell activation – EC50 (nM) | 526 |
| Mouse A2A T cell activation – EC50 (nM) | 229 |
| **Permeability for oral therapy** | |
| Caco-2 A->B (10-6 cm/sec.) (efflux ratio) | 6.4 (1.7) |
| **Lipophilicity** | |
| LogD (pH 7.5) | 1.5 |
| **Metabolic stability** | |
| Human Microsomes Clint, app (µL.min/mg) | <12 |
| Human Hepatocytes Clint, app (µL.min/mg/10-6 cells) | 4 |

### Unmet Need

A2A receptor blockade has the potential to target a multitude of tumour types including colorectal, NSCLC, renal cell cancer, or RCC, triple-negative breast cancer, TNBC, and many others. In patients with recurrent and/or metastatic solid tumours treated with immune checkpoint inhibitors, or ICIs, (i.e., anti-PD1/PDL1/CTL4) only 25% have durable responses. In 2018, there were 1.8 million incident cases of NSCLC worldwide, and the number is expected to increase to 1.9 million by 2027. ICIs are approved for NSCLC, but a significant number of patients fail to respond, and there are few options available for patients who progress.

### Our Approach

We set out to design a potent, highly selective antagonist with low CNS penetration. Several unique elements of our AI platform enabled us to achieve this goal:
- our SPR biosensor expertise enabled us to run a fragment screen on wild-type receptors for the target, generating novel chemical equity;
- we rapidly performed 2D evolutions on these fragments to generate highly potent and selective molecules;
- our platform incorporated knowledge from the published 3D structures to refine selectivity;
- we further confirmed selectivity by developing a suite of new assays for a broad range of adenosine receptors, including $A_{2A}$, $A_{2B}$, $A_1$, $A_3$, CD73 and CD39, which generated key insights across four pre-clinical species.

### Our Solution

Our AI-first approach generated a highly differentiated $A_{2A}$ antagonist with notable activity against the target, high selectivity, low-CNS penetration, and high tumour exposure (see Figure 3 below). Our eventual clinical candidate, EXS21546, was identified within nine months of generating novel designs, and we identified our candidate after testing just 163 compounds. Figure 1 below shows our drug candidate is highly selective for $A_{2A}$ receptors while also demonstrating other favourable design attributes. This high selectivity translated into straightforward pharmacology with a saturable concentration response in functional assays, compared to competitor molecules (see Figure 4). In addition, EXS21546 was shown to be selective over a large panel of GPCRs, ion channels, transporters and kinases. Our drug candidate exhibited the desired PK profile with low CNS penetration (see Figure 3) when compared to some competing approaches. In a pre-clinical study, EXS21546 demonstrated comparable single agent anti-tumour activity to an approved anti-PD-1 (see Figure 2 below). The positive pre-clinical data for this drug candidate illustrates the ability of our AI approach to rapidly find a potential solution to a difficult treatment challenge. In December 2020, we initiated our Phase 1 clinical trial for EXS21546.
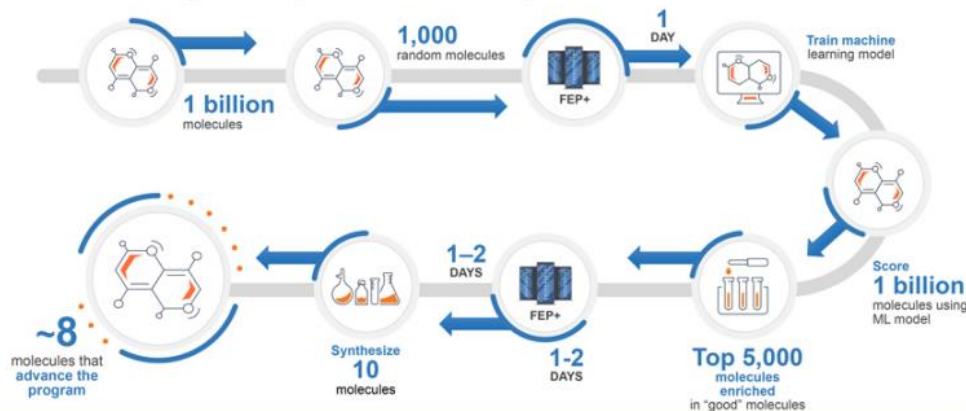
Figure source: Exscientia reports

# Schrodinger
## – AI/ML Clinical Candidate
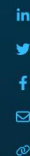


Schrödinger platform combines accuracy of physics with speed of machine learning

Enables ultra-large scale exploration of chemical space

1 billion molecules → 1,000 random molecules → FEP+ → 1 DAY → Train machine learning model → Score 1 billion molecules using ML model → Top 5,000 molecules enriched in "good" molecules → 1-2 DAYS → FEP+ → 1-2 DAYS → Synthesize 10 molecules → ~8 molecules that advance the program

SCHRÖDINGER.



FIERCE Biotech

Research   Biotech   Medtech   CRO   Special Reports   Trending Topics   Podcasts

MEDTECH

**Schrödinger to kick off first human trial of its computer-designed blood cancer drug**

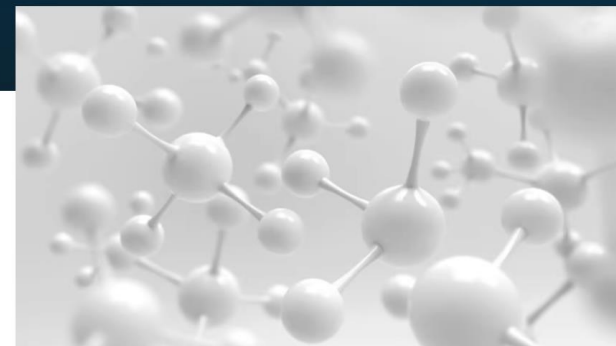By Conor Hale • Jun 29, 2022 02:16am

Schrodinger   Artificial Intelligence   drug discovery   clinical research

Is it a machine learning platform developer, a clinical biotech, or both simultaneously? Schrödinger has received FDA permission to move forward with its first in-human study. ((Getty Images))

After working with biotechs and Big Pharmas to accelerate their research programs, high-tech molecule modeler Schrödinger is taking its first steps as a clinical company itself. It has secured an FDA green light to study its computer-designed therapy for non-Hodgkin

Figure source: Schrodinger.com

# Schrodinger
## – AI/ML Clinical Candidate

**8.2 billion**
compounds computationally evaluated

**78**
total compounds synthesized in lead series

**10 months**
to discovery of development candidate

" The ability to leverage the computational platform to rapidly identify not just one, but several novel, highly potent series with well-balanced properties is unique in my many years experience in industry."

—**Zhe Nie,** Project Lead
Executive Director, Medicinal Chemistry,
Schrödinger Therapeutics Group

| | |
|---|---|
| **Target** | MALT1, protease |
| **Program Type** | Schrödinger proprietary program, small molecule |
| **Indication** | Relapsed or refractory B-cell lymphoma, chronic lymphocytic leukemia |
| **Stage** | Phase 1 clinical trial |

MEDTECH

## Schrödinger to kick off first human trial of its computer-designed blood cancer drug

By **Conor Hale** • Jun 29, 2022 02:16am

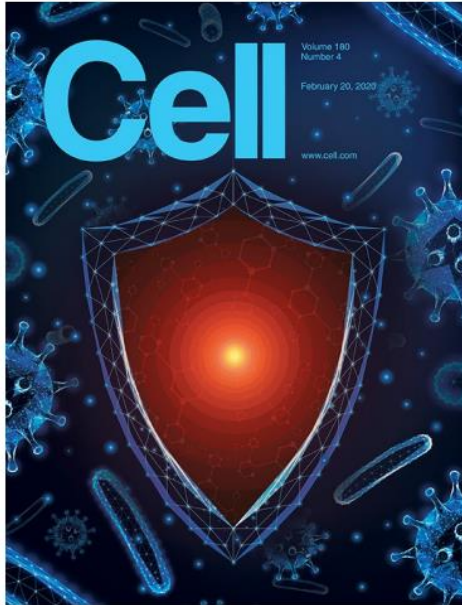Schrodinger    Artificial Intelligence    drug discovery    clinical research

Is it a machine learning platform developer, a clinical biotech, or both simultaneously? Schrödinger has received FDA permission to move forward with its first in-human study. ((Getty Images))

After working with biotechs and Big Pharmas to accelerate their research programs, high-tech molecule modeler Schrödinger is taking its first steps as a clinical company itself. It has secured an FDA green light to study its computer-designed therapy for non-Hodgkin

**NLP SUMMIT**

Figure source: Schrodinger.com

# And Many More Success Stories



Figure Source: Stokes, Yang, Swanson, Jin et al, Cell 2020

# Drug Discovery – Complex, Challenging Process

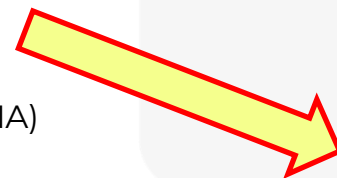Search in the chemical space

A good drug (e.g., kills virus)

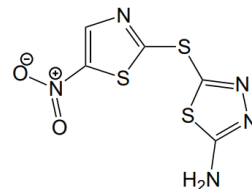**Biological target**: macromolecule (such as., protein, RNA) involved in the disease pathway
**Compound/Hit/Lead/Candidate** : Binds and interact with the target
**Toxicity** : isn't harmful to organism
**Selectivity** : binding specifically to desired biological target and many more such as., **Solubility, Caco2, ADME, t1/2, and PK/PD**
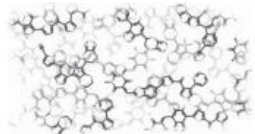
# Title

**Chemical space**



(Drug-like, photovoltaics, polymers, dyes)

Search in the chemical space

**Functional space**

Desired properties (redox potential, solubility, toxicity)

A good drug (e.g., kills virus)
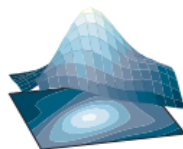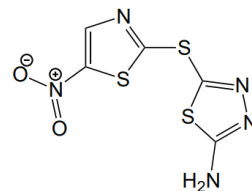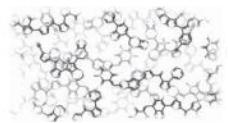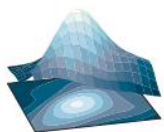
# Drug Discovery – Complex, Challenging Process

**Chemical space**

(Drug-like, photovoltaics, polymers, dyes)

**Functional space**

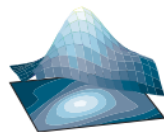Desired properties (redox potential, solubility, toxicity)

**Functional space**

Desired properties (redox potential, solubility, toxicity)

**Chemical space**

(Drug-like, photovoltaics, polymers, dyes)

NLP SUMMIT

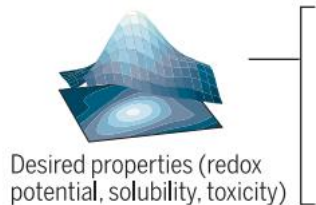# Drug Designing – Strategies

# Molecular Representations



Molecular Formula
$C_{13}H_{18}O_2$

List of fragments
C(=O)O
c1ccccc1
CC(C)C
...

Electrostatics

Molecule

Molecular graph

3D coordinates (x, y, z)

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

# Molecular Representations



N1CCN(CC1)C(C(F)=C2)=CC(=C2C4=O)N(C3CC3)C=C4C(=O)O

SMILES string for ciprofloxacin

- SMILES (Simplified molecular-input line-entry system)
- 

  string representation

# Generative AI methods in DD



Drug Discovery Today: Technologies

# Generative AI methods in DD



Fig. The relative frequencies of the various AI/ML model architectures observed in the review and their chronological progression.

# Title

# REINVENT – WORKFLOWS for denovo design

# REINVENT – WORKFLOWS for denovo design



- The generative model is subjected to transfer learning with a smaller set of compounds that are relevant to the project of interest
- This will bias the resulting model to produce project specific compounds with much higher probability than any random compounds
- Therefore much smaller dataset can suffice to find good hits when using the scoring function

# REINVENT – WORKFLOWS for denovo design



- Transfer learning is conducted with a set of compounds that have the desired properties
- For example, if we aim to maximize a predictive model among the other components we would use all the compounds that are considered as active by this model; If we aim towards certain subseries of compounds we would only use those that share the specific features for transfer learning
- After concluding with transfer learning the resulting agent is **"focused"** on the specific set
- We can conduct TL for multiple epochs & observe the stats from each epoch thus deciding which agent is focused enough

# **Benchmarking?**

| Property | Practical relevance | Rule-based | Distribution-based |
|---|---|---|---|
| *Validity*—molecules must adhere to chemical principles e.g., valency. | Critical. | Molecules should always be valid (unless there are systematic errors in hard-coded rules). | Dependent on molecular representation chosen, complexity of training data, and complexity of model. |
| *Uniqueness*—the rate at which molecules are duplicated by the model. | Unnecessary if the single de novo molecule satisfies all desirable properties. | Dependent on the search algorithm used. | Dependent on the search algorithm and applicability domain imposed by training data. |
| *Diversity*—the scope of chemotypes generated relative to all chemical space. | Unnecessary if de novo molecules occupy the most optimal chemical space. | Dependent on the search algorithm and fidelity achievable by chemical building rules (i.e., atoms or fragments etc.). | Dependent on the search algorithm and applicability domain imposed by training data. May afford greater diversity where rules are difficult to explicitly define (e.g., natural products). |
| *Novelty*—the presence of molecules in any training data used. | Critical to fulfill the definition of de novo molecule generation. | Only applicable to seeded models such as genetic algorithms. | Dependent on all model aspects, training data used, molecular representation, architecture, etc. |
| *Similarity*—the similarity between generated molecules and any training data. | Unnecessary if de novo molecules satisfy all desirable properties. | Only applicable to seeded models such as genetic algorithms. | Dependent on all model aspects, training data used, molecular representation, architecture, etc. |
| *Synthetic feasibility*—the ability to synthesize a molecule in the lab with relative ease. | Critical for experimental validation and practical application as a therapeutic. | Rules can adhere to known chemical reactions and reactive sites, ensuring a degree of synthetic feasibility (usually to the detriment of diversity). | Synthetic feasibility of molecules may be implicitly learned based on the training data; however, it cannot be guaranteed for novel molecules. |

| Target | Training | Approach | Cpds Generated | Filters | Screening? | Active Hits |

**2017**

VEGFR2

**University of Oxford**

VEGFR2 actives (25,000) from Binding DB

char-RNN

10,000 Compounds

• Solubility
• Synthesizability
• Similarity with VEGFR2 actives
• Docking Score Cutoff

5 Compounds Tested

2 Compounds Active @ (1-1000 nM)

**2018**

PPARs & RXRs agonists

**ETH Zurich**

25 RXR/PPAR fatty acid Active Mimetics

LSTM

1000 Compounds (−COOH)

• SPIDER – target prediction
• BB's availability
• Similarity with known actives
• Docking Score Cut-off

5 Compounds Tested

4 Compounds Active @ (0.06 - 14 uM)

**2018**

Kinase Panel

**Insilico Medicine**

30,000 Kinase inhibitors from Thomson Integrity DB

ATNC

30,000 Compounds

• MedChem Filters
• Clustering
• Reward function for diversity
• Similarity with known actives (> 0.7)

50 cpds Tested

7 cpds Active @ (100% inh 10 uM)

NLP SUMMIT

| Target | Training | Approach | Cpds Generated | Filters | Screening? | Active Hits |
|--------|----------|----------|----------------|---------|------------|-------------|

**2018**

JAK2/JAK3

**Insilico Medicine**

Known actives → **SSAAE** → 30,000 Compounds →

- Docking Score
- Off target effects
- Property Analysis
- Molecular dynamics
- MedChem Feedback

→ 1 Compound Tested

1 Compound Active @ (6 M) →



**2018**

DDR1

**Insilico Medicine**

17,000 Kinase inhibitors + DDR1 active cpds → **GENTRL** → 30,000 Compounds →

- MedChem Filters
- Clustering
- Pharmacophore
- Synthesizability
- MedChem filters

→ 6 cpds Tested

4 cpds Active @ nM range →



NLP SUMMIT

# Limitations and future scope

- **Generation in Low-data Regime**

- **Lack of Unified Evaluation Protocols**

- **Lack of Large-scale Study and Benchmark**

- **Out-of-distribution Generation**

- **Unrealistic Problem Formulation**

- **Ideal assumptions/expensive Oracle Calls**

- **Lack of Interpretability**