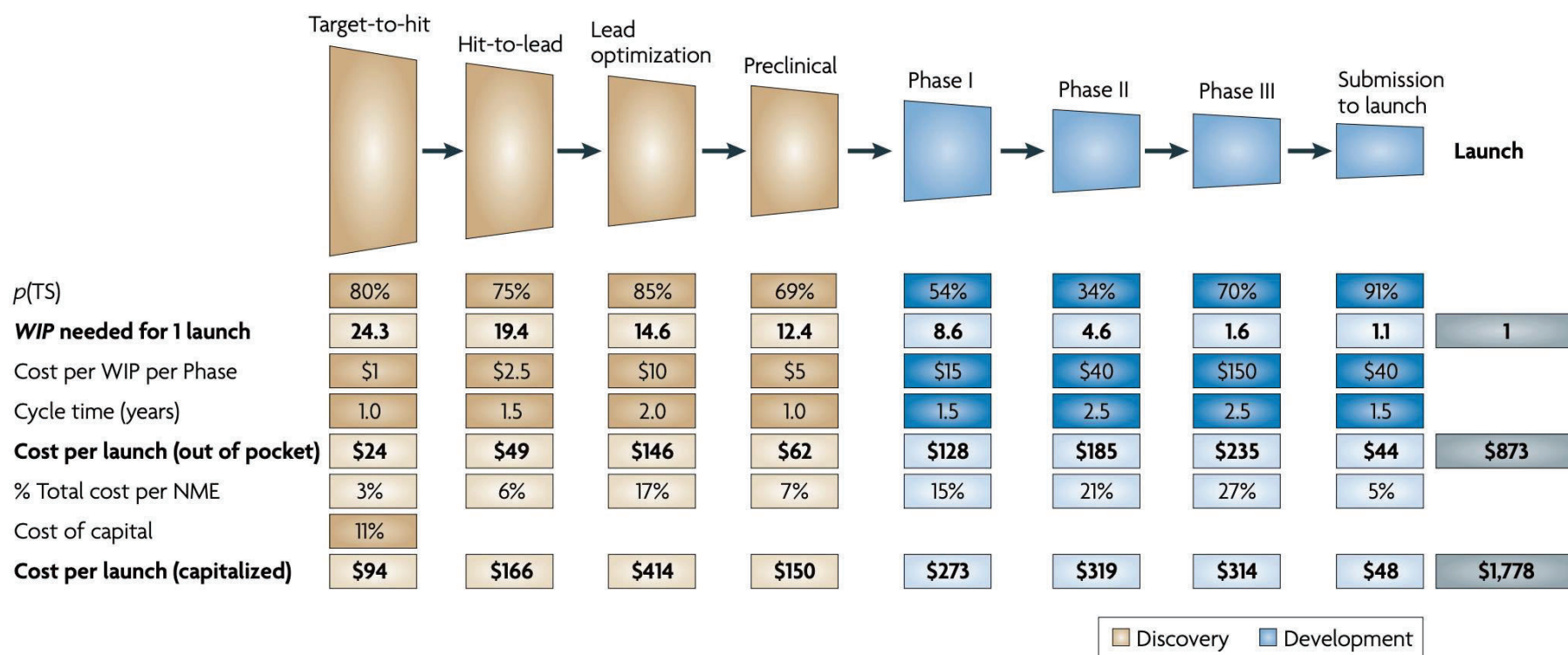


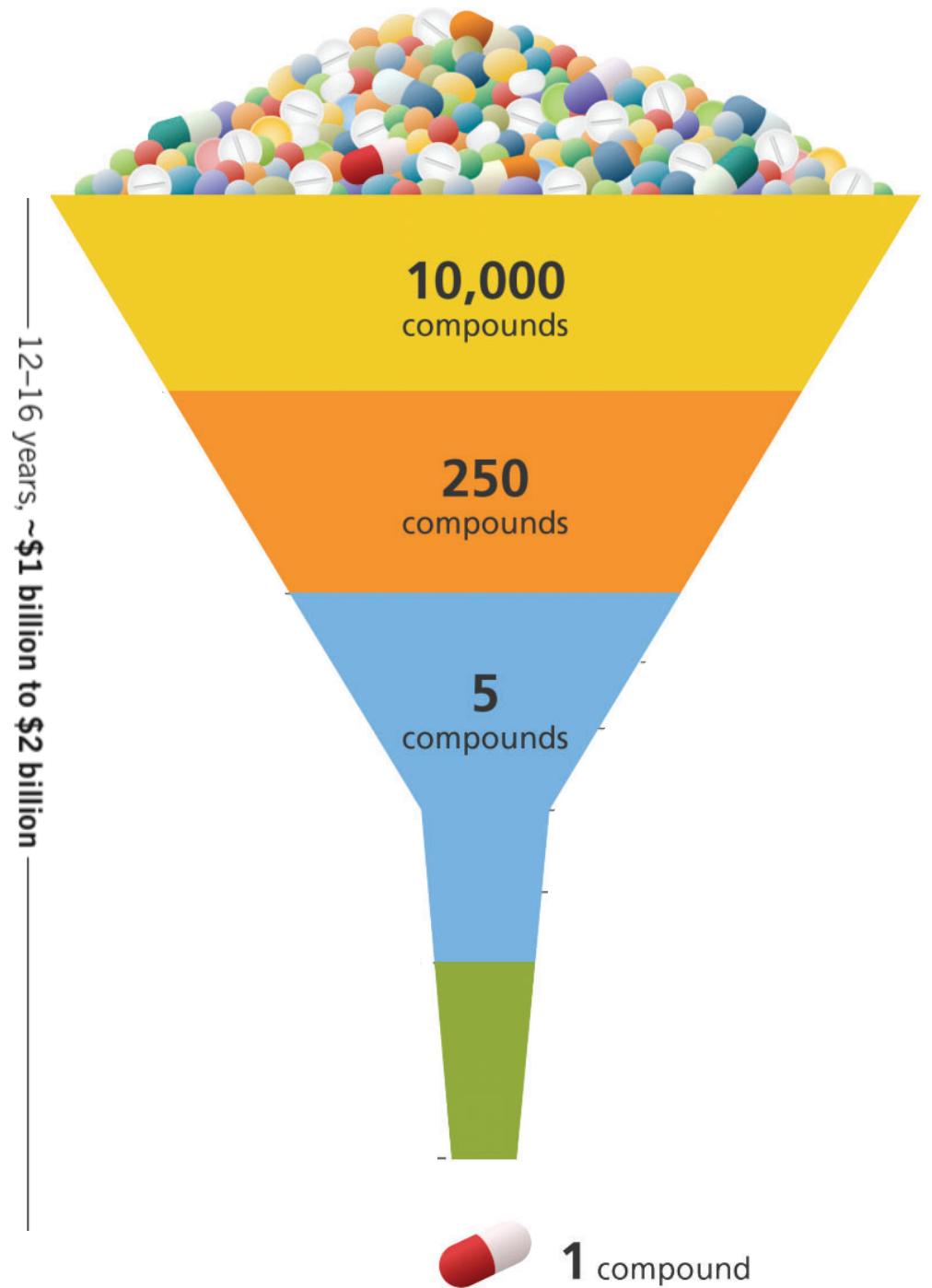
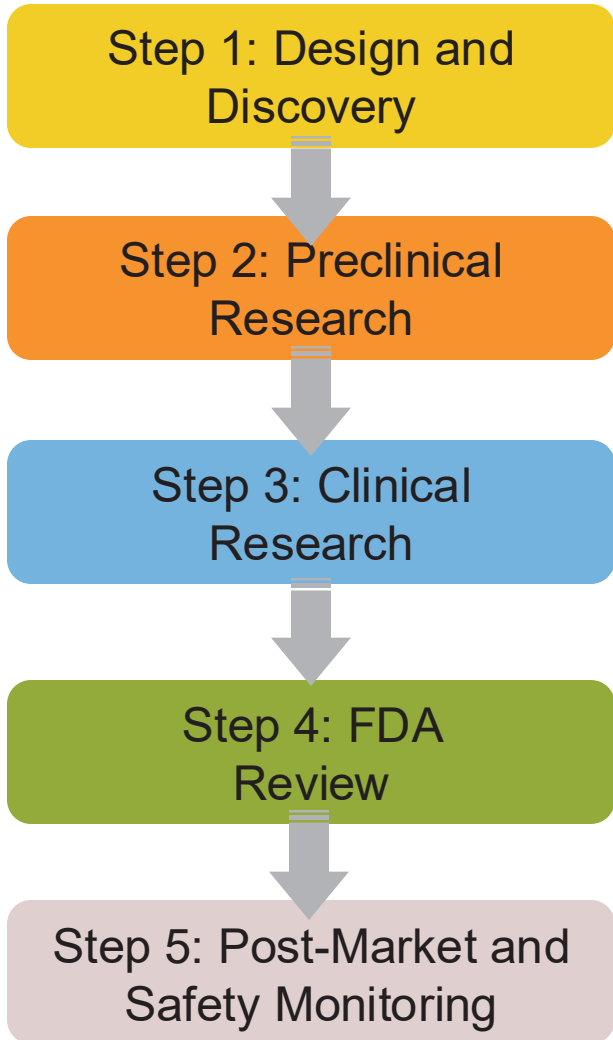
Phases of drug discovery from initial stage (target-to-hit) to final stage (launch)

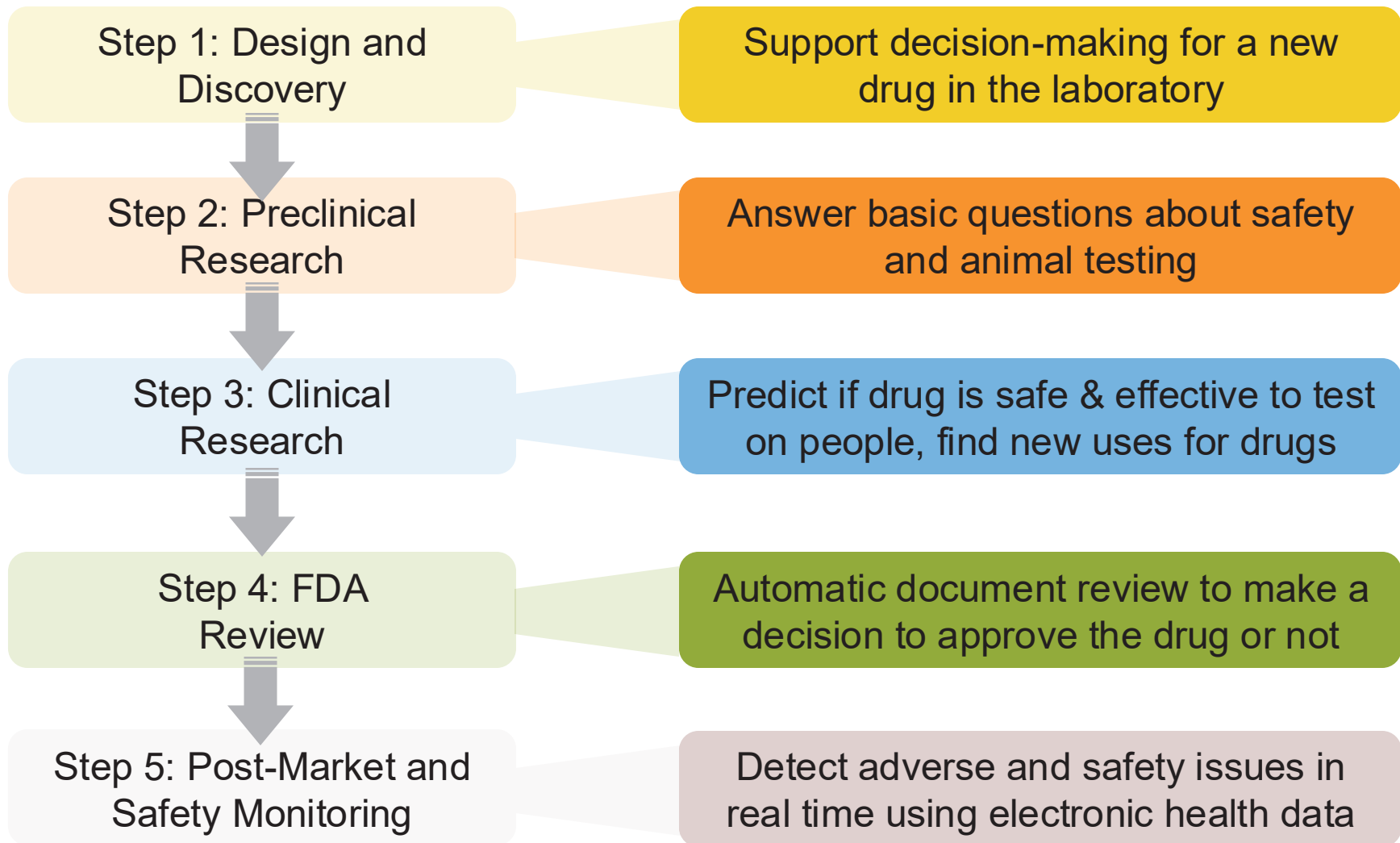


$p(\text{TS})$ – probability of successful transition from one stage to the next; NME – new molecular entity; WIP – work in process

Drug-like chemical space
 10^{60} chemical compounds

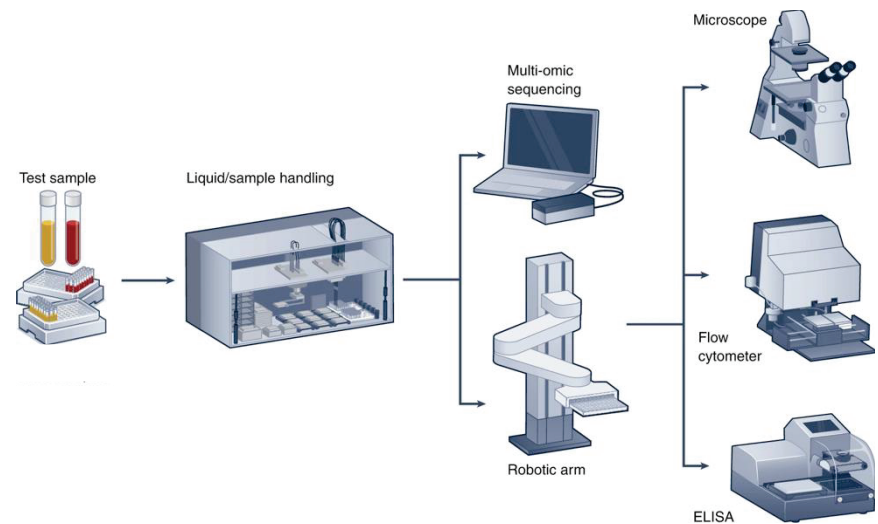
**Drugs available
to humans**
 $\sim 10^4$ •





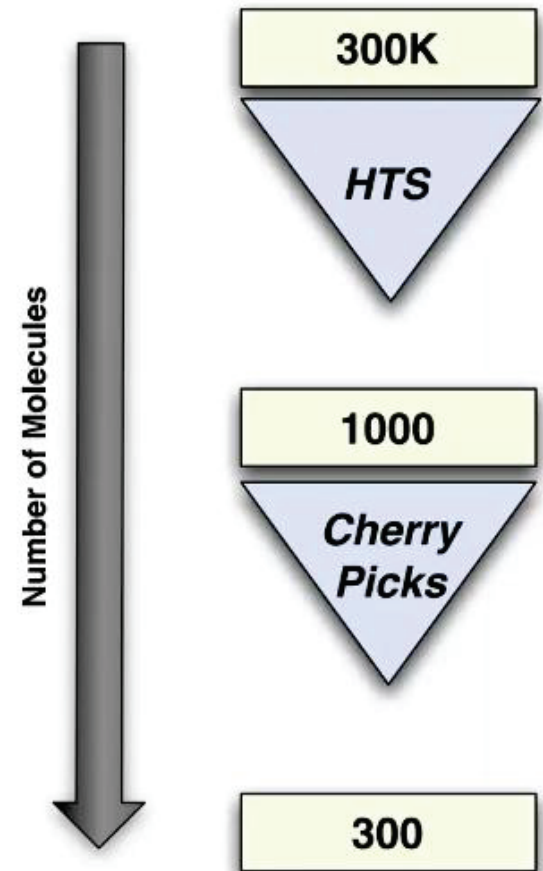
High throughput screening (HTS)

- Test thousands to hundreds of thousands of compounds in one or more assays
 - Biochemical, genetic, and pharmacological assays
- Integrate with robotics for self-driving lab
- **Goal:** Rapidly identify novel modulators of biological systems
 - Cellular basis of diseases
 - Therapeutic agents



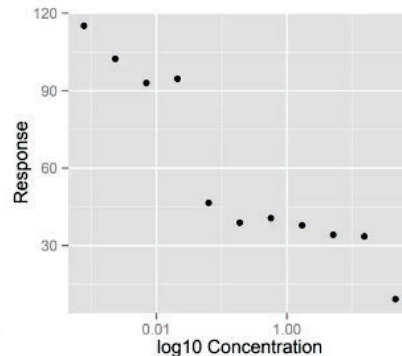
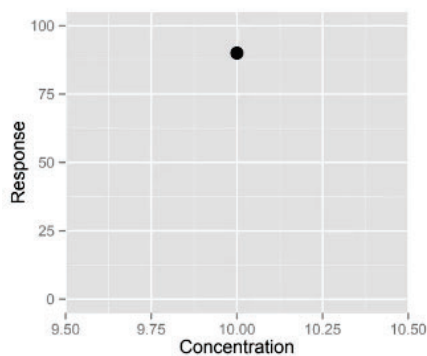
Goals of high throughput screening

- Rapidly screen large collections of compounds (chemical libraries)
- Efficiently identify active compounds
 - Test them in slower, accurate, expensive screens
- Use the data to learn what types of compounds tend to be active

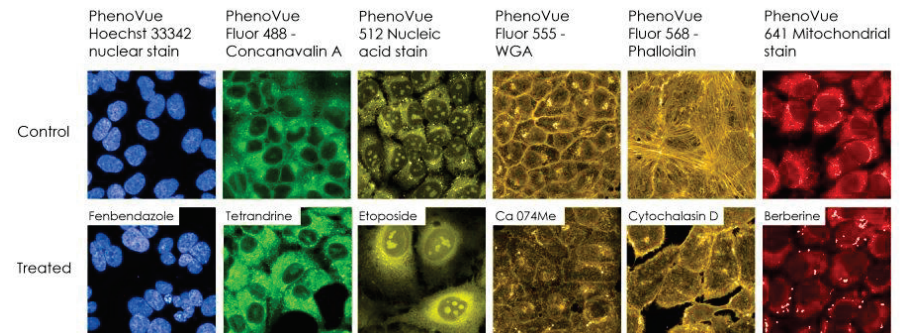


HTS data types

- Categorical: active/inactive or toxic/nontoxic
- Continuous: single-point or dose-response
- Multiple readouts:
 - Might read at different wavelengths or time points
 - More complex when dealing with images



Single-point vs. dose-response readouts



Cell painting for phenotypic drug discovery

HTS: Machine learning setup

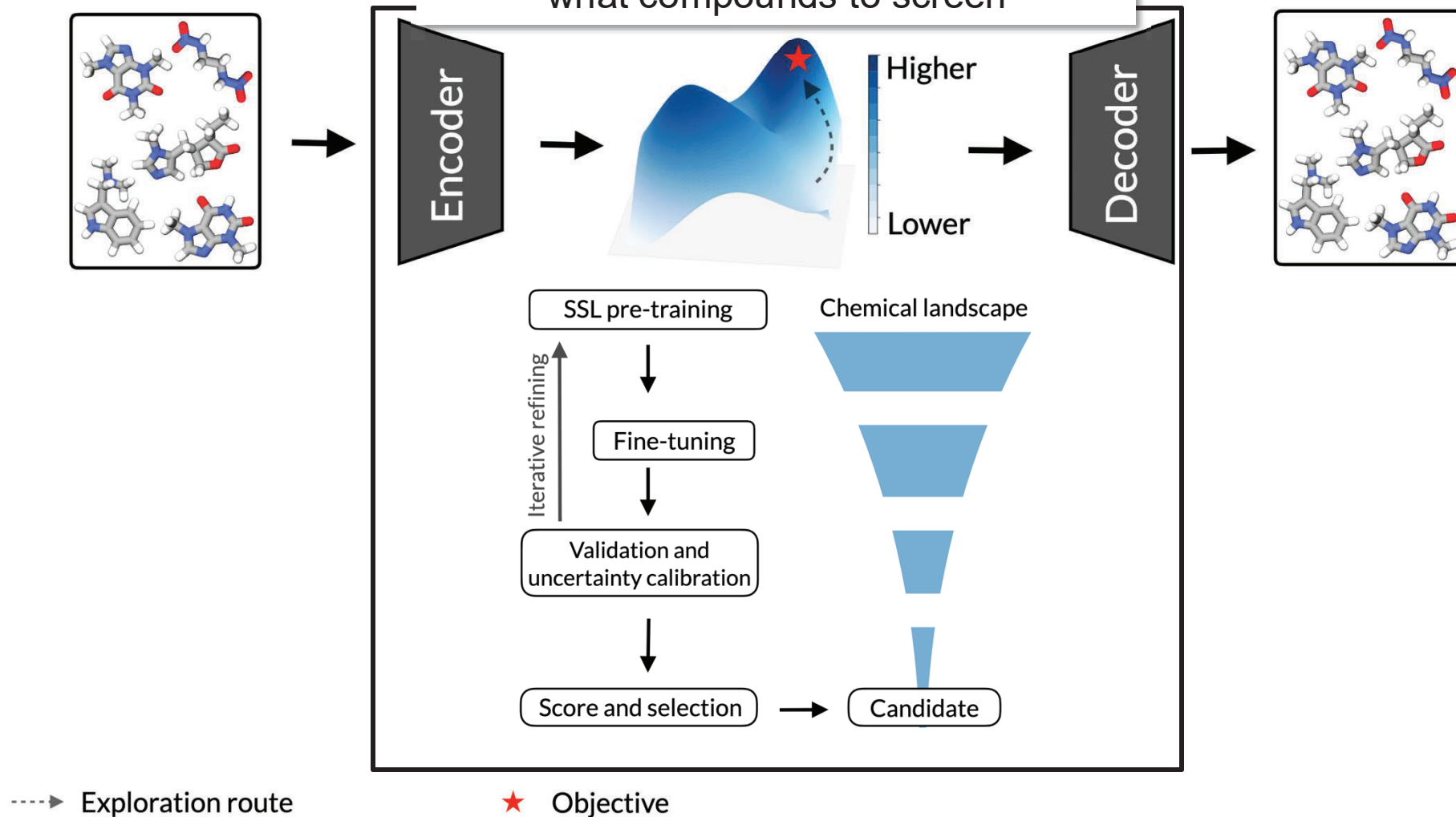
- HTS tests the activity of molecules:

$$\textit{Activity} = f(\textit{Structure})$$

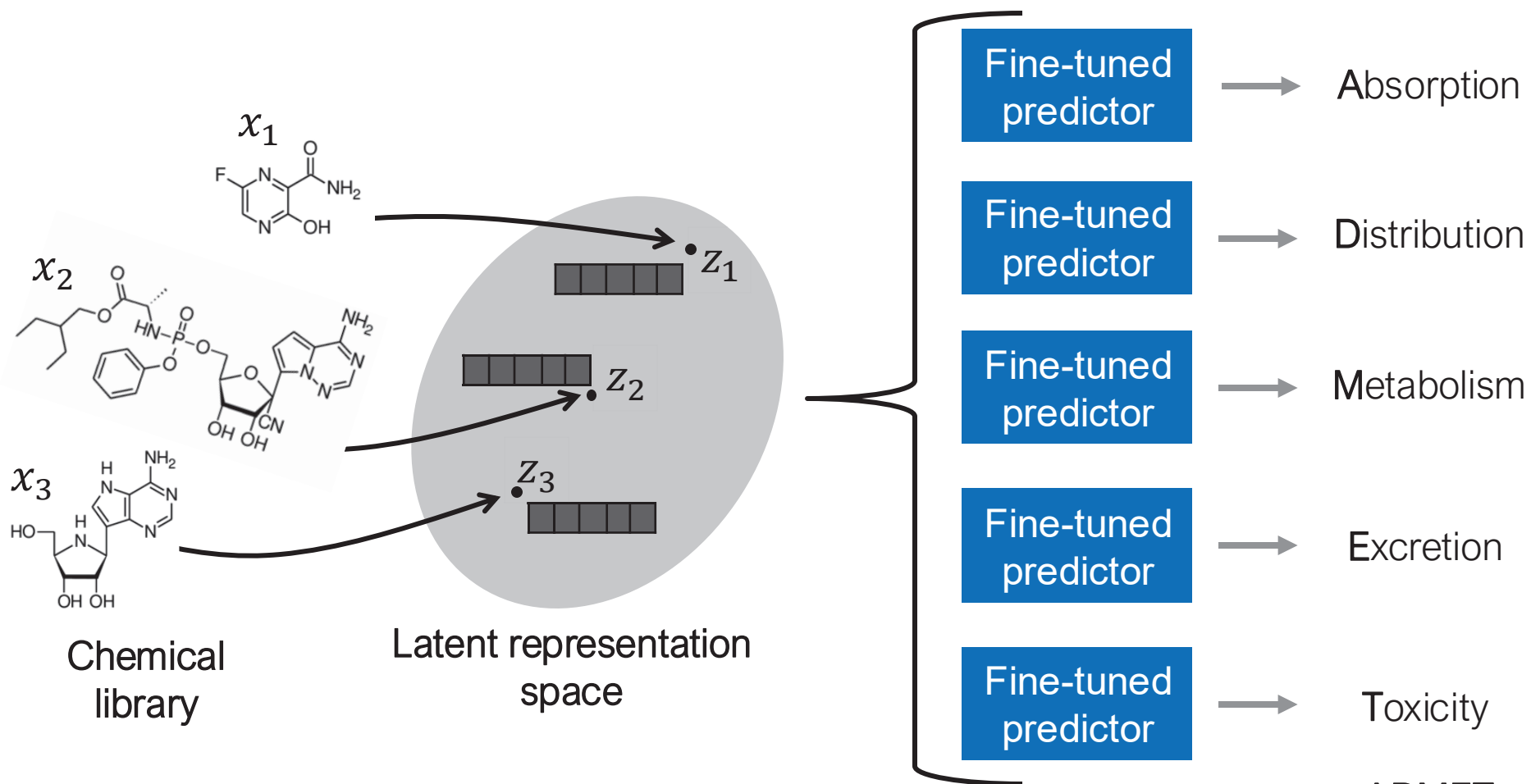
- We need to describe the molecular structure
 - Various discrete or real-valued descriptors
 - Surfaces (3D)
 - Binary fingerprints
 - Learned molecular embeddings

In-silico screening and optimization of molecular structure

Use computational models to suggest what compounds to screen



Molecular property prediction



Self-supervised learning: Molecules with similar molecular structure get embedded close together. Various representations: Neural fingerprints, Attentive fingerprints, SMILES descriptors, Graphormer, Transformer-M, and others

ADMET endpoints

What can we use molecular representations for?

- **Search**

- Given a potent active molecule, find similar ones (or dissimilar but also potent)

- **Prediction of various endpoints**

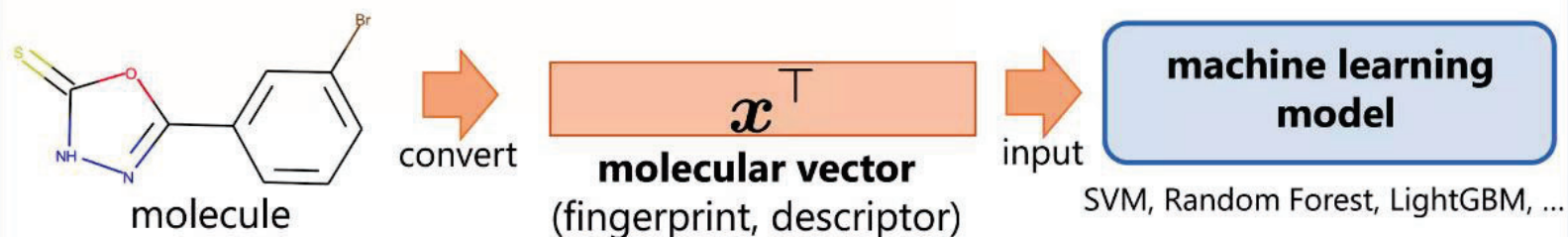
- Given a set of active and inactive molecules, build a model to predict which members from a chemical library will be active

- **Clustering**

- Given a set of molecules, do they cluster into structurally different groups?

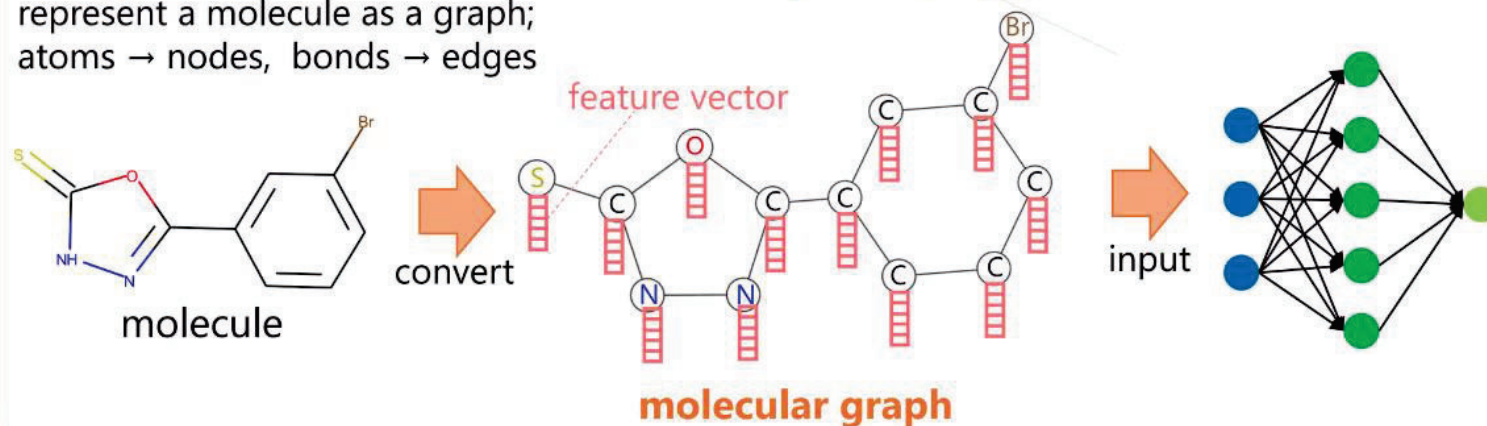
Two strategies for producing molecular representations

Traditional approach

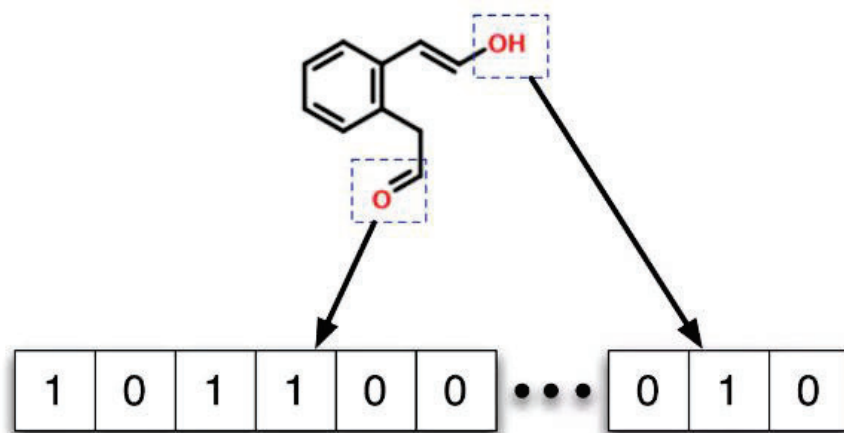


Graph convolutional network (GCN) approach

represent a molecule as a graph;
atoms → nodes, bonds → edges



Fingerprint representations



- Lots of types of fingerprints
- Keyed fingerprints indicate the presence or absence of a structural feature
- Length can vary from 166 to 4096 bits or more
- Fingerprints usually compared to each other using the Tanimoto metric

Towards neural fingerprints

Algorithm 1 Circular fingerprints

```
1: Input: molecule, radius  $R$ , fingerprint length  $S$ 
2: Initialize: fingerprint vector  $\mathbf{f} \leftarrow \mathbf{0}_S$ 
3: for each atom  $a$  in molecule
4:    $\mathbf{r}_a \leftarrow g(a)$   $\triangleright$  lookup atom features
5: for  $L = 1$  to  $R$   $\triangleright$  for each layer
6:   for each atom  $a$  in molecule
7:      $\mathbf{r}_1 \dots \mathbf{r}_N = \text{neighbors}(a)$ 
8:      $\mathbf{v} \leftarrow [\mathbf{r}_a, \mathbf{r}_1, \dots, \mathbf{r}_N]$   $\triangleright$  concatenate
9:      $\mathbf{r}_a \leftarrow \text{hash}(\mathbf{v})$   $\triangleright$  hash function
10:     $i \leftarrow \text{mod}(r_a, S)$   $\triangleright$  convert to index
11:     $\mathbf{f}_i \leftarrow 1$   $\triangleright$  Write 1 at index
12: Return: binary vector  $\mathbf{f}$ 
```

Algorithm 2 Neural graph fingerprints

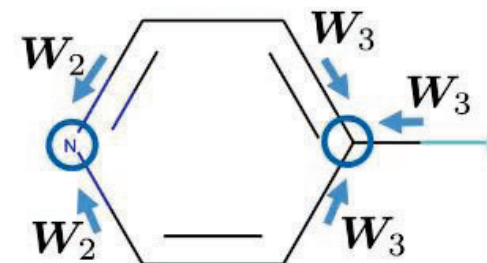
```
1: Input: molecule, radius  $R$ , hidden weights  $H_1^1 \dots H_R^5$ , output weights  $W_1 \dots W_R$ 
2: Initialize: fingerprint vector  $\mathbf{f} \leftarrow \mathbf{0}_S$ 
3: for each atom  $a$  in molecule
4:    $\mathbf{r}_a \leftarrow g(a)$   $\triangleright$  lookup atom features
5: for  $L = 1$  to  $R$   $\triangleright$  for each layer
6:   for each atom  $a$  in molecule
7:      $\mathbf{r}_1 \dots \mathbf{r}_N = \text{neighbors}(a)$ 
8:      $\mathbf{v} \leftarrow \mathbf{r}_a + \sum_{i=1}^N \mathbf{r}_i$   $\triangleright$  sum
9:      $\mathbf{r}_a \leftarrow \sigma(\mathbf{v} H_L^N)$   $\triangleright$  smooth function
10:     $\mathbf{i} \leftarrow \text{softmax}(\mathbf{r}_a W_L)$   $\triangleright$  sparsify
11:     $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{i}$   $\triangleright$  add to fingerprint
12: Return: real-valued vector  $\mathbf{f}$ 
```

Figure 2: Pseudocode of circular fingerprints (*left*) and neural graph fingerprints (*right*). Differences are highlighted in blue. Every non-differentiable operation is replaced with a differentiable analog.

Neural fingerprint representations

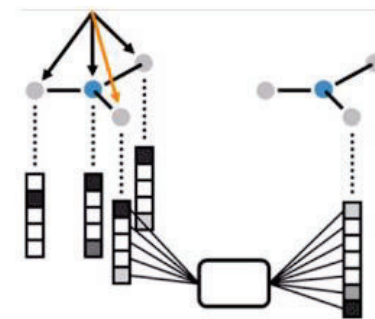
1) Neural graph fingerprints

- Generate molecular fingerprints with a neural network
- Update atom features using only adjacent atoms
- Use different weights for node degrees



2) Molecular graphs

- Update atom features by convolutional and pooling layers using adjacent atoms



- They did not consider property of edges (bonds)
- They did not consider atoms other than 1-neighbor