

1 **Supplementary Information – Multi-Constraint Molecular**  
2 **Generation using Sparsely Labelled Training Data for**  
3 **Localized High-Concentration Electrolyte Diluent Screening**

4 **Jonathan P. Mailoa,<sup>1\*</sup> Xin Li,<sup>1</sup> Jiezhong Qiu,<sup>1</sup> and Shengyu Zhang<sup>2\*</sup>**

5

6 1) Tencent Quantum Laboratory, Tencent, Shenzhen, Guangdong, China

7 2) Tencent Quantum Laboratory, Tencent, Hong Kong SAR, China

8

9 \* corresponding author: [jpmailoa@alum.mit.edu](mailto:jpmailoa@alum.mit.edu), [shengyuzhang@tencent.com](mailto:shengyuzhang@tencent.com)

10

## 1 Baseline SSVAE Model Cost Function Derivation

2 Our work mostly utilizes the cost function originally developed by Kang, *et al* based on the work  
3 by Kingma, *et al*. The regression loss component  $R_{SSVAE}(x,y)$  from Equation 3 in the main text is a  
4 straightforward regression square loss function and will not be discussed in more detail.

5 The VAE cost function for the labelled entries (**Equation 1**) originates from **Equation 1** in Kang  
6 paper and **Equation 6** in Kingma paper. It is described that the variational lower bound  $-L(x,y)$  of the  
7 log-probability of a labelled instance  $(x,y)$  is:

$$\begin{aligned} \ln p(x,y) &\geq E_{q_\phi(z|x,y)}[\ln p_\theta(x|y,z) + \ln p(y) + \ln p(z) - \ln q_\phi(z|x,y)] \\ &= E_{q_\phi(z|x,y)}[\ln p_\theta(x|y,z)] + \ln p(y) - D_{KL}(q_\phi(z|x,y)||p(z)) \\ &=-L(x,y) \\ 8 \quad L(x,y) &=-E_{q_\phi(z|x,y)}[\ln p_\theta(x|y,z)] - \ln p(y) + D_{KL}(q_\phi(z|x,y)||p(z)) \end{aligned}$$

9 The first term is simply the cross-entropy loss. When summed over all the  $n_L$  fully labelled samples in  
10 the minibatch and the  $n_x$  dimension corresponding to our molecule SMILES one-hot vector  
11 representation (see Kang paper for detail), the first term is converted to **(1)**:

$$12 \quad -\sum_{i=1}^{n_L} \sum_{j=1}^{n_x} (x_{i,j} \ln x_{D,i,j} + (1 - x_{i,j}) \ln (1 - x_{D,i,j}))$$

13 For the second term, we first remember that  $p(y)$  is a multivariate normal distribution with mean  
14 matrix  $E$  ( $n_y \times 1$ ) and covariance matrix  $C$  ( $n_y \times n_y$ ) constructed from all the available training labels in  
15 the dataset, where  $n_y$  is the dimension of  $y$ . The multivariate probability distribution function is:

$$16 \quad p(y) = \frac{1}{\sqrt{(2\pi)^{n_y} \det(C)}} e^{-\frac{1}{2}(y-E)^T C^{-1} (y-E)}$$

17 When the term  $(-\ln p(y))$  is expanded for all molecule labels  $y_L$  in the minibatch, we have **(2)**:

$$18 \quad \sum_{i=1}^{n_L} \frac{1}{2} \left( n_y \ln 2\pi + \ln(\det(C)) + \sum_{j=1}^{n_y} (y_{L,i,j} - E_j) \sum_{k=1}^{n_y} (y_{L,i,k} - E_k) C_{k,j}^{-1} \right)$$

19

1 For the third term, we recall that the definition of Kullback-Leibler divergence loss is:

2 
$$D_{KL}(P||Q) = \sum_{x \in \mathcal{X}} P(x) \ln \left( \frac{P(x)}{Q(x)} \right)$$

3 We also remember that the prior distribution  $p(z) = N(z|0,I)$  and approximated posterior distribution

4  $q_\phi(z|x,y) = N(z|\mu_\phi(x,y), \text{diag}(\sigma_\phi^2(x,y)))$  are used in the Kang paper. Before we proceed with the

5 derivation, recall the Gaussian function integral formulas  $\int_{-\infty}^{\infty} e^{-\alpha x^2} dx = \sqrt{\pi/\alpha}$ ,  $\int_{-\infty}^{\infty} x e^{-\alpha x^2} dx = 0$  and

6  $\int_{-\infty}^{\infty} x^2 e^{-\alpha x^2} dx = \frac{1}{2} \sqrt{\pi/\alpha^3}$ . Also recall that  $\int_{-\infty}^{\infty} e^{-\alpha(x-\beta)^2} dx = \sqrt{\pi/\alpha}$ , and  $\int_{-\infty}^{\infty} (x-\beta) e^{-\alpha(x-\beta)^2} dx = 0$ , and

7  $\int_{-\infty}^{\infty} (x-\beta)^2 e^{-\alpha(x-\beta)^2} dx = \frac{1}{2} \sqrt{\pi/\alpha^3}$  hold for the constants  $\alpha$  and  $\beta$ . For the sake of clarity, in the following

8 derivation we will drop  $x, y$ , and the subscript  $\phi$  from the functions and just use  $z$ , with  $n_z$  being its

9 dimension. Furthermore, note that the prior and approximated posterior distributions are

10 independent between different  $z$  dimensions, so we can process them individually.  $i$  and  $j$  are the

11 labelled molecule and the  $z$  dimension indexes, and we will drop this subscript in the integral for

12 clarity:

$$q(z_{i,j}) = \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(z_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}}$$

13 
$$p(z_j) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z_j^2}{2}}$$

$$\begin{aligned}
D_{KL}(q(z_{i,j})||p(z_j)) &= \int_{-\infty}^{\infty} \frac{1}{\sigma_{i,j}\sqrt{2\pi}} e^{-\frac{(z-\mu_{i,j})^2}{2\sigma_{i,j}^2}} \ln \left( \frac{1}{\sigma_{i,j}} e^{-\frac{z^2}{2} - \frac{(z-\mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) dz \\
&= \frac{1}{\sigma_{i,j}\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(z-\mu_{i,j})^2}{2\sigma_{i,j}^2}} \left( -\ln \sigma_{i,j} + \frac{(z-\mu_{i,j})^2}{2} + (z-\mu_{i,j})\mu_{i,j} + \frac{\mu_{i,j}^2}{2} - \frac{(z-\mu_{i,j})^2}{2\sigma_{i,j}^2} \right) dz \\
&= -\ln \sigma_{i,j} + \frac{\sigma_{i,j}^2}{2} + 0 + \frac{\mu_{i,j}^2}{2} - \frac{1}{2}
\end{aligned}$$

1

The readers can

2 convince themselves that repeating the exercise above using the joint distribution  $D_{KL}(q(z_i)||p(z))$  will  
3 simply generate the summation of the term over  $j$  (the  $n_z$  dimensions of  $z$ ). After we consider the  
4 summation over  $i$  and  $j$  and put back the  $z$  notation on  $\mu$  and  $\sigma$  for clarity, we get the following term  
5 **(3)**, completing **Equation 1** of the main text:

$$- \sum_{i=1}^{n_L} \sum_{j=1}^{n_z} \frac{1}{2} (1 + \ln \sigma(z_{i,j})^2 - \mu(z_{i,j})^2 - \sigma(z_{i,j})^2)$$

6

7 The VAE cost function for the unlabelled molecule entries (**Equation 2**) originates from **Equation**  
8 **2** in Kang paper and **Equation 7** in Kingma paper. It is described that the variational lower bound  
9  $-U(x,y)$  of the log-probability of an unlabelled instance ( $x$ ) is:

$$\begin{aligned}
\ln p(x) &\geq E_{q_{\phi}(y,z|x)} [\ln p_{\theta}(x|y,z) + \ln p(y) + \ln p(z) - \ln q_{\phi}(y,z|x)] \\
&= E_{q_{\phi}(y,z|x)} [\ln p_{\theta}(x|y,z) + \ln p(y) - \ln q_{\phi}(y|x) + \ln p(z) - \ln q_{\phi}(z|x,y)] \\
&= E_{q_{\phi}(y,z|x)} [\ln p_{\theta}(x|y,z)] - D_{KL}(q_{\phi}(y|x)||p(y)) - E_{q_{\phi}(y|x)} [D_{KL}(q_{\phi}(z|x,y)||p(z))] \\
&= -U(x) \\
U(x) &= -E_{q_{\phi}(y,z|x)} [\ln p_{\theta}(x|y,z)] + D_{KL}(q_{\phi}(y|x)||p(y)) + E_{q_{\phi}(y|x)} [D_{KL}(q_{\phi}(z|x,y)||p(z))]
\end{aligned}$$

10

11 The derivation for the first and the third term are identical to their counterparts for the fully  
12 labelled molecules. They will correspondingly generate **(4)**:

$$- \sum_{i=1}^{n_U} \sum_{j=1}^{n_x} (x_{i,j} \ln x_{D,i,j} + (1-x_{i,j}) \ln (1-x_{D,i,j}))$$

13

14 and **(5)**:

$$1 \quad - \sum_{i=1}^{n_U} \sum_{j=1}^{n_z} \frac{1}{2} (1 + \ln \sigma(z_{i,j})^2 - \mu(z_{i,j})^2 - \sigma(z_{i,j})^2)$$

2           The second loss function term related to unlabelled  $\mathcal{Y}$  is significantly more complicated. Recall  
3 that we have the multivariate probability prior distribution function for  $\mathcal{Y}$  based on the fully labelled  
4 molecules:

$$5 \quad p(\mathcal{Y}) = \frac{1}{\sqrt{(2\pi)^{n_y} \det(C)}} e^{-\frac{1}{2}(\mathcal{Y}-E)^T C^{-1}(\mathcal{Y}-E)}$$

6 Also recall that we have taken the assumption of normal posterior distribution for the sampled  $\mathcal{Y}$   
7 generated by the predictor  $q_\phi(\mathcal{Y}|x) = N(\mathcal{Y}|\mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$  below.  $i$  and  $j$  correspond to unlabelled  
8 molecule and  $\mathcal{Y}$  dimension indices, and we have dropped  $\mathcal{Y}$  from  $\mu$  and  $\sigma$  for clarity:

$$9 \quad q(y_{i,j}) = \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}}$$

10 The complex form for the joint distribution of  $p(\mathcal{Y})$  means we no longer have the luxury of separating  
11 the integrals based on the label space dimension the way we have done for  $p(z)$ .

$$12 \quad D_{KL}(q(y_i)||p(\mathcal{Y})) = \int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \ln \left( \frac{\sqrt{\det(C)}}{\prod_j \sigma_{i,j}} e^{\left( \frac{1}{2}(\mathcal{Y}-E)^T C^{-1}(\mathcal{Y}-E) - \sum_j \frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2} \right)} \right) dy_i$$

$$= \int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \left( \ln \frac{\sqrt{\det(C)}}{\prod_j \sigma_{i,j}} + \frac{1}{2}(\mathcal{Y}-E)^T C^{-1}(\mathcal{Y}-E) - \sum_{j=1}^{n_y} \frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2} \right) dy_i$$

13 Note that this integration is performed over the  $n_y$  dimensions of  $\mathcal{Y}_i$ . We separate the terms for clarity.

14 The first term in the parenthesis is simply a constant, and integrating over it generates **(6)**:

$$15 \quad \int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \ln \frac{\sqrt{\det(C)}}{\prod_j \sigma_{i,j}} dy_i = \frac{1}{2} \ln(\det(C)) - \frac{1}{2} \sum_{j=1}^{n_y} \ln \sigma_{i,j}^2$$

- 1 The third term in the parenthesis consists of  $n_y$  terms. For each of them (let's say  $j = k$ ), the following  
 2 situation applies:

$$\begin{aligned}
 & - \int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2} dy_i \\
 & = - \int_{-\infty}^{\infty} \left( \prod_{j \neq k} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) dy_{i,j \neq k} \int_{-\infty}^{\infty} \left( \frac{1}{\sigma_{i,k} \sqrt{2\pi}} e^{-\frac{(y_{i,k} - \mu_{i,k})^2}{2\sigma_{i,k}^2}} \right) dy_{i,k} = - (1^{n_y - 1}) \frac{1}{2}
 \end{aligned}$$

3

- 4 After summation over  $j$ , we simply get  $-\frac{n_y}{2}$ , which is the summation term **(7)**.

- 5 Finally, for the third term in the parenthesis, we recognize that it is simply a double summation over  
 6  $n_y^2$  terms.

$$\begin{aligned}
 & (y - E)^T C^{-1} (y - E) \\
 & = \sum_{j=1}^{n_y} \sum_{k=1}^{n_y} C_{kj}^{-1} (y_{i,j} - E_j) (y_{i,k} - E_k) = \sum_{j=1}^{n_y} \sum_{k=1}^{n_y} C_{kj}^{-1} ((y_{i,j} - \mu_{i,j})(y_{i,k} - \mu_{i,k}) + (y_{i,j} - \mu_{i,j})(\mu_{i,k} - E_k) + (\mu_{i,j} - E_j)(y_{i,k} - \mu_{i,k}) + (\mu_{i,j} - E_j)(\mu_{i,k} - E_k))
 \end{aligned}$$

7

- 8 Based on the Gaussian function integral formula, remember that the terms  $(y_{i,j} - \mu_{i,j})(\mu_{i,k} - E_k)$  and  
 9  $(\mu_{i,j} - E_j)(y_{i,k} - \mu_{i,k})$  will contribute 0 to the integral, so we will drop them.

- 10 When  $j = k$ , we have the integral term **(8)**:

$$\int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \left( \frac{1}{2} C_{jj}^{-1} ((y_{i,j} - \mu_{i,j})^2 + (\mu_{i,j} - E_j)^2) \right) dy_i = \frac{1}{2} C_{jj}^{-1} \sigma_{i,j}^2 + \frac{1}{2} C_{jj}^{-1} (\mu_{i,j} - E_j)^2$$

11

- 12 When  $j \neq k$ , the first term will additionally generate zero contribution to the integral, and we simply  
 13 have the integral term **(9)**:

$$\int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \left( \frac{1}{2} C_{kj}^{-1} ((\mu_{i,j} - E_j)(\mu_{i,k} - E_k)) \right) dy_i = \frac{1}{2} C_{kj}^{-1} (\mu_{i,j} - E_j)(\mu_{i,k} - E_k)$$

14

- 1 Now that we have all the individual components within the double summation integrated, we can  
 2 perform double summation of **(8)** and **(9)** to get the integral term **(10)**:

$$3 \int_{-\infty}^{\infty} \left( \prod_{j=1}^{n_y} \frac{1}{\sigma_{i,j} \sqrt{2\pi}} e^{-\frac{(y_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}} \right) \left( \frac{1}{2} (\mathbf{y} - \mathbf{E})^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{E}) \right) d\mathbf{y}_i = \frac{1}{2} \sum_{j=1}^{n_y} C_{jj}^{-1} \sigma_{i,j}^2 + \frac{1}{2} \sum_{j=1}^{n_y} (\mu_{i,j} - E_j) \sum_{k=1}^{n_y} (\mu_{i,k} - E_k) C_{k,j}^{-1}$$

- 4 Now that we have all the integral component terms **(6)**, **(7)**, **(10)** for  $\mathcal{Y}$ , after summation over all the  
 5 unlabelled molecules  $n_U$  and putting back the  $\mathcal{Y}_{P,i,j}$  into  $\mu$  and  $\sigma$  for completeness, we have the  
 6 following term **(11)**:

$$7 \sum_{i=1}^{n_U} \frac{1}{2} \left( \sum_{j=1}^{n_y} C_{jj}^{-1} \sigma(\mathcal{Y}_{P,i,j})^2 + \sum_{j=1}^{n_y} (\mu(\mathcal{Y}_{P,i,j}) - E_j) \sum_{k=1}^{n_y} (\mu(\mathcal{Y}_{P,i,k}) - E_k) C_{k,j}^{-1} - n_y + \ln(\det(\mathbf{C})) - \sum_{j=1}^{n_y} \ln \sigma(\mathcal{Y}_{P,i,j})^2 \right)$$

- 8 This Kullback-Leibler divergence loss term for  $\mathcal{Y}$  completes our re-derivation of the original baseline  
 9 SSSVAE model cost function as implemented by Kang, *et al*, with **(4)**, **(11)**, and **(5)** forming the total VAE  
 10 loss for the unlabelled molecules in **Equation 2** of the main text.

11

12

13

14

## 1 Impact of ConGen Missing Label Imputation

2

3 We use the same training dataset which we have used in the main text **Table 2**, containing two  
4 databases: 1) ZINC database with 310k molecules containing  $Mol.Wt$ ,  $LogP$ , and  $QED$ , and 2) Materials  
5 Project electrolyte molecule database containing  $Mol.Wt$ ,  $IE$ , and  $EA$ . We trained two different RNN-  
6 based ConGen models, one utilizing the mean and covariance table ( $E$  and  $C$ ) imputation (the one used  
7 in the main text), and a second model where such imputation throughout the training is not  
8 performed. Afterward, we perform three different queries: **(1)**  $Mol.Wt = 250 Da$ ,  $LogP = 2.5$ ,  $IE = 5.0$   
9  $eV$ , **(2)**  $LogP = 2.5$ ,  $IE = 5.0 eV$ , and **(3)**  $Mol.Wt = 250 Da$ ,  $LogP = 2.5$ ,  $EA = 4.0 eV$ . For each query task,  
10 we obtain 1000 molecules (3-5% of generated molecules are invalid) and validate the generated  
11 molecules'  $Mol.Wt$  and  $LogP$  using RDKit.  $IE$  and  $EA$  are not validated due to the computational costs,  
12 but these conditions are used for the conditional generation tasks in this section because the impact  
13 of imputation only manifests when there are insufficient valid samples for  $C$  calculation (such as  
14  $LogP-IE$  pairs). From **Supplementary Table 1** below, we can see that imputation has relatively no  
15 impact on task **(1)**, positive impact on task **(2)**, and negative impact on task **(3)**. It is expected that  
16 imputation helps on task **(2)**, as there is almost no entry overlap to calculate the correlation matrix  
17 entry between  $LogP$  and  $IE$  (the only two query constraints), making the  $C$  matrix with imputation  
18 statistically more meaningful. However, we are also able to find a situation where imputation is  
19 harmful, such as task **(3)** where the query condition  $EA = 4.0 eV$  is near the extreme end of the property  
20 label distribution. Future study should be conducted to determine the circumstances where  
21 imputation usage will be beneficial more systematically.

Query	RNN Model Type	$Mol.Wt (Da)$	$LogP$
<b>(1)</b> $Mol.Wt = 250$ , $LogP = 2.5$ , $IE = 5.0$	Imputation	$250 \pm 4$	$2.54 \pm 0.31$
	No imputation	$250 \pm 4$	$2.53 \pm 0.31$
<b>(2)</b> $LogP = 2.5$ , $IE = 5.0$	<b>Imputation</b>	$314 \pm 90$	<b><math>2.54 \pm 0.40</math></b>
	No imputation	$318 \pm 94$	$2.44 \pm 0.43$
<b>(3)</b> $Mol.Wt = 250$ , $LogP = 2.5$ , $EA = 4.0$	Imputation	$251 \pm 8$	$2.97 \pm 0.53$
	<b>No imputation</b>	$249 \pm 8$	<b><math>2.78 \pm 0.50</math></b>

22 **Supplementary Table 1 | Impact of missing label imputation on different RNN ConGen multi-constraint**  
23 **generation queries.** Best model for each task query is bolded. For task **(1)**, the two models are equivalent.



## 1 Impact of ConGen BERT $\beta$ Variations and Training from Scratch

2 We use the same training dataset which we have used in the main text **Table 2**, to train BERT-  
3 type ConGen model (transferred ChemBERTa parameters used just like in the main text, with just the  
4 last layer being unfrozen). In the main text, we use  $\beta = 10000$  for both RNN and BERT-based ConGen,  
5 the same setting used to train the original RNN-based SSVAE model but fail to see clear model  
6 improvement when BERT-based ConGen is used. Correspondingly, we no longer use BERT-based  
7 ConGen in the subsequent sections as it is computationally more expensive than RNN-based ConGen.  
8 In this supplementary section, we attempt different  $\beta$  settings (1000, 3000, 10000, 30000, and  
9 100000) to see if we can get transfer learning BERT-based ConGen to clearly perform better than RNN-  
10 based ConGen. We also train a BERT-based ConGen from scratch with the original  $\beta = 10000$  to see its  
11 impact on performance. From **Supplementary Table 2** below, we can see that  $\beta = 10000$  seems to  
12 have the best ConGen BERT performance, while training the ConGen BERT from scratch instead of  
13 utilizing transfer learning produces a model with poor regression and conditional generation  
14 capabilities.

Model	$\beta$	<i>Mol.Wt</i> (Da)	<i>LogP</i>	<i>QED</i>	<i>EA</i> (eV)	<i>IE</i> (eV)
<b>Predictor MAE</b>						
RNN	10000	2.70	0.05	0.009	0.20	0.16
BERT	1000	7.09	0.16	0.020	0.25	0.24
	3000	6.24	0.15	0.017	0.22	0.19
	10000	6.07	0.15	0.017	0.22	0.19
	10000 – fresh start	82.6	0.95	0.129	0.85	1.17
	30000	6.33	0.16	0.017	0.23	0.20
	100000	6.17	0.15	0.017	0.22	0.19
<b>Conditional Generation</b>						
RNN	10000	248 ± 4	2.55 ± 0.23	0.672 ± 0.082	2.06 ± 0.55	6.53 ± 0.62
BERT	1000	229 ± 36	2.55 ± 0.61	0.716 ± 0.071	2.26 ± 0.93	6.35 ± 0.45
	3000	209 ± 34	2.75 ± 0.64	0.751 ± 0.029	1.36 ± 0.50	6.48 ± 0.24
	10000	252 ± 3	2.45 ± 0.36	0.756 ± 0.127	1.80 ± 0.64	6.36 ± 0.41
	10000 – fresh start	249 ± 6	2.45 ± 0.64	0.659 ± 0.125	1.20 ± 1.15	6.64 ± 0.47
	30000	253 ± 9	3.01 ± 0.26	0.790 ± 0.075	1.62 ± 0.75	6.42 ± 0.36
	100000	252 ± 5	2.44 ± 0.46	0.668 ± 0.131	2.39 ± 0.98	6.54 ± 0.74

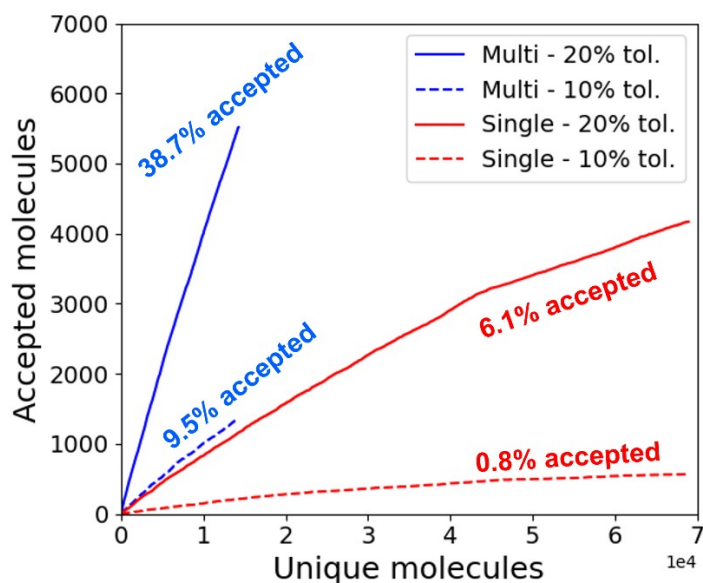
15 **Supplementary Table 2 | Impact of  $\beta$  variation on BERT-based ConGen performance.** We use the same  
16 conditional molecule generation with multiple constraints we use in **Table 2** of the main text: *Mol.Wt* = 250 Da,  
17 *LogP* = 2.5, and *IE* = 5.e eV.



## 1 ConGen Advantage on Multi-Condition Generative Design

2 It is instructive to contemplate whether enabling multi-condition generative model such as  
3 ConGen is useful or not, compared to single-condition generative models. We use the RNN-based  
4 ConGen model we have trained and utilized in the main text **Table 2**, with the intent of generating  
5 molecules which have the following simultaneous properties:  $Mol.Wt = 250$ ,  $LogP = 2.5$ , and  $QED =$   
6  $0.55$ . In the multi-constraint approach, we query the model 90,000 times with these simultaneous  
7 constraints. Out of these 90,000 molecules, 4,257 are within the training dataset and we generate a  
8 total of 14,628 unique molecules outside of the training dataset. In the single-constraint approach, we  
9 query the model 30,000 times for each of the single constraint. Out of these 90,000 molecules, 3,928  
10 are within the training dataset and we generate a total of 68,925 unique molecules outside of the  
11 training dataset. Note that multi-constraint ConGen here generates less unique molecules compared  
12 to single-constraint version, but the generated molecule properties are more accurate as we will show  
13 below. Because it is impossible for the model to generate molecule with the same exact properties as  
14 our reference values, we need to decide on an acceptance tolerance criterion. If our tolerance  
15 criterion is 20% relative error, we obtain 5,518 and 4,170 acceptable molecules from the multi-  
16 constraint and single-constraint approach respectively, corresponding to 38.7% and 6.1% acceptance  
17 rate for the two approaches. If our tolerance criterion is 10% relative error (stricter), we only obtain  
18 1,354 and 568 acceptable molecules from the multi-constraint and single-constraint approach  
19 respectively, corresponding to 9.5% and 0.8% acceptance rate for the two approaches. It is  
20 straightforward to see the user's tolerance on the generated molecule properties determine which  
21 approach should be used. If the user has very large tolerance (in extreme case, any molecule property  
22 error is accepted), then the single-constraint approach should be used because the model generates  
23 more diverse and unique molecules. However, if the user has strict requirements on the generated  
24 molecule properties, a multi-constraint approach will generate acceptable molecules with significantly  
25 higher acceptance rate. We demonstrate this result in **Supplementary Figure 1** below. It is important  
26 to note that the 10% and 20% relative error acceptance criterion we use here is simply for convenience

- 1 purposes. In general, the user should specify the acceptance criterion which makes more sense for
- 2 each molecule property in the user application.



3

4 **Supplementary Figure 1 | Acceptance rate of multi-constraint ConGen vs single-constraint ConGen.** The same  
5 generative model is used for both approaches, but the multi-constraint version is queried 90000 times using the  
6 3 simultaneous constraints while the single-constraint version is queried 30000 times for each of the 3 single  
7 constraints. The single-constraint approach can generate more diverse and unique molecules because it is less  
8 restricted, but the acceptance rate of the generated molecule properties is much lower than the multi-constraint  
9 generative model approach. Both acceptance under the 20% tolerance and the stricter 10% tolerance criteria  
10 are shown.

11

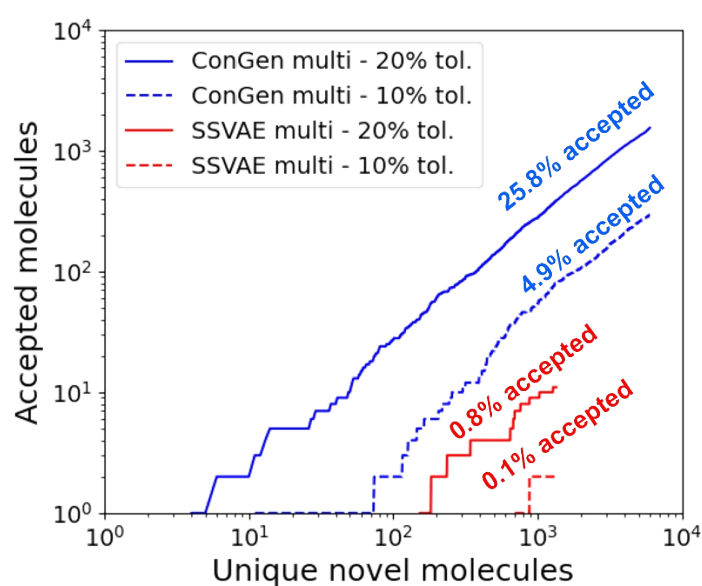
## 1 ConGen Advantage Compared to SSVAE on Incomplete-Labelled Dataset

2 In the main text, we have mentioned that the primary advantage of ConGen is that it can work  
3 with molecule databases with incomplete labels, which is not suitable for a baseline SSVAE model.  
4 Suppose we would like to utilize the SSVAE model for multi-condition generative modelling tasks using  
5 training molecule databases with incomplete labels regardless, for comparison purposes. We can  
6 perform this task with two different approaches on the SSVAE model before comparing its  
7 performance to a ConGen model trained on the same datasets:

- 8 1. Use all molecules with full labels as fully labelled training dataset and designate any  
9 molecules with incomplete labels as fully unlabelled training dataset.
- 10 2. Train individual SSVAE models with only single-property label each, ensuring that each  
11 SSVAE model can utilize all the training property labels.

12 In the first approach, we allow the SSVAE model to perform multi-constraint molecule  
13 generation, in exchange for a significant loss of property training data label information. It is  
14 impossible to test this approach on the main text's training datasets because of the extremely low  
15 availability of fully labelled molecules in our dataset. We have instead utilized the full ZINC database  
16 from the main text (containing 310,000 unique molecules) with fully labelled *Mol.Wt*, *LogP*, and *QED*  
17 properties (obtained using RDKit). 300,000 of these molecules are designated as the training dataset,  
18 while the remaining 10,000 molecules are designated as the test dataset. For each property label  
19 column in the training dataset, we randomly de-label 70% of the properties. Consequently, only ~2.7%  
20 of the molecules in our training dataset are fully labelled (8,117 fully labelled + 291,883 fully unlabelled  
21 molecules). This reflects the severe consequence of randomness we typically encounter from available  
22 experimental databases when multiple properties are needed. We train the SSVAE model of Kang, *et*  
23 *al* on this dataset, and perform a query to generate 10,000 molecules in multi-constraint mode (  
24 *Mol.Wt* = 250, *LogP* = 2.5, and *QED* = 0.55). Note that the original SSVAE model as published by Kang,  
25 *et al* only supports single-constraint molecule generation mode, and we have slightly modified the

1 SSSVAE model's molecule generation function to enable multi-constraint generation mode the exact  
2 same way it is being done in ConGen. We also train the ConGen model on the same partially de-  
3 labelled training dataset (ConGen can utilize all ~70% remaining partial labels) and perform a query to  
4 generate 10,000 molecules with same multi-property constraints as specified above. Out of the 10,000  
5 SSSVAE molecules, 488 molecules are within the training dataset and we obtain 1,347 unique molecules  
6 outside of the training dataset. On the other hand, ConGen generates 239 molecules which are within  
7 the training dataset and 5,988 unique molecules outside of the training dataset. If our tolerance  
8 criterion is 20% relative error, we obtain 11 and 1,545 acceptable molecules from the SSSVAE and  
9 ConGen approaches respectively, corresponding to 0.8% and 25.8% acceptance rate for the two  
10 approaches. If our tolerance criterion is 10% relative error (stricter), we only obtain 2 and 295  
11 acceptable molecules from the SSSVAE and ConGen approaches respectively, corresponding to 0.1%  
12 and 4.9% acceptance rate for the two approaches. See Supplementary **Figure 2** below for more details.



13

14 **Supplementary Figure 2 | Acceptance rate of multi-constraint ConGen vs multi-constraint SSSVAE.** The same  
15 training dataset (ZINC database with *Mol.Wt*, *LogP*, and *QED* labels) is used for both approaches, and a query  
16 for generating 10,000 molecules with multi-property constraints (*Mol.Wt* = 250 Da, *LogP* = 2.5, and *QED* =  
17 0.55) is performed using both models. However, the multi-constraint SSSVAE cannot utilize labels from molecules  
18 which are only partially labelled. Because of the lack of fully labelled molecules, the multi-constraint SSSVAE  
19 model suffers on both its ability to generate diverse molecules under multi-property constraints, and on its  
20 ability to generate molecules with the correct desired multi-property constraints. Both acceptance under the  
21 20% tolerance and the stricter 10% tolerance criteria are shown.

1 In the second approach, we restrict the individual SSSVAE model training to be done on only one  
 2 property each. In this example, we can utilize the training dataset we have previously used in the main  
 3 text **Table 2** and perform direct comparison with the molecules generated by the ConGen approach.  
 4 We train three individual SSSVAE models for each of the three property constraints (*Mol.Wt* = 250 Da,  
 5 *LogP* = 2.5, and *IE* = 5 eV), and then generate 10 molecules for each SSSVAE model's validation using  
 6 the corresponding single-property constraint. Note that the original SSSVAE model as published by  
 7 Kang, *et al* will generate programming errors when trained on single properties (the code as published  
 8 was trained on multiple properties and can only be used for single-property generation tasks), and  
 9 some minor code reprogramming is needed to enable the model to work on single-property training  
 10 and molecule generation tasks. The result is show in **Supplementary Table 3**, where we show that the  
 11 individual SSSVAE models separately trained on *Mol.Wt* and *IE* have good control over the individual  
 12 property of the molecule it generates, but surprisingly the model separately trained on the single  
 13 *LogP* property has bad performance on *LogP* molecule generation task (*LogP* =  $2.35 \pm 1.12$ ). While  
 14 the individual SSSVAE models have good performance on the single property constraint it was trained  
 15 on, they are not constrained on the other two properties they are not trained on. Consequently, these  
 16 molecules generated by these single-constraint SSSVAE's are not suitable for satisfying the requirement  
 17 of multi-constraint generation queries, compared to the ConGen model.

RNN-based Model	<i>Mol.Wt</i> (Da)	<i>LogP</i>	<i>QED</i>	<i>EA</i> (eV)	<i>IE</i> (eV)
<b>ConGen</b>	<b>248 ± 4</b>	<b>2.55 ± 0.23</b>	0.672 ± 0.082	2.06 ± 0.55	<b>6.53 ± 0.62</b>
SSVAE ( <i>MolWt</i> = 250 Da)	<b>251 ± 7</b>	2.68 ± 0.82	0.780 ± 0.091	1.76 ± 0.62	6.14 ± 1.39
SSVAE ( <i>LogP</i> = 2.5)	300 ± 82	<b>2.35 ± 1.12</b>	0.745 ± 0.149	2.04 ± 0.79	5.94 ± 0.82
SSVAE ( <i>IE</i> = 5.0 eV)	322 ± 57	2.81 ± 1.07	0.716 ± 0.154	2.54 ± 0.78	<b>6.27 ± 0.50</b>
<b>SSVAE (combined)</b>	<b>291 ± 65</b>	<b>2.61 ± 1.03</b>	0.747 ± 0.137	2.11 ± 0.78	<b>6.12 ± 0.95</b>

18 **Supplementary Table 3 | ConGen comparison with baseline SSSVAE models trained using single properties on**  
 19 **multi-constraint decoder conditional generation task.** We use the same conditional molecule generation with  
 20 multiple constraints we use in **Table 2** of the main text: *Mol.Wt* = 250 Da, *LogP* = 2.5, and *IE* = 5.e eV. The  
 21 ConGen model is trained on these 3 property types simultaneously, while the SSSVAE models are trained on  
 22 individual property types. Note that while the generated molecules' *IE* magnitudes do not perfectly match the  
 23 *IE* constraint (likely due to the difference in quantum chemistry workflow between ours and the online  
 24 database), ConGen and SSSVAE (*IE*) models both demonstrate tight *IE* control for the generated molecules.

# 1 Complete List of Molecules Generated in Main Text Query 1 (Figure 4b)

2 Supplementary Table 4 | Candidate Li-ion battery LHCE diluent molecules generated with multi-constraint  
 3 ConGen model (Query 1).

['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 250, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 250, 4.0, 1.0]
Cc1ccc(OCC(F)(F)C(F)(F)Cl)cc1 Nc1c(F)cc(OC(F)(F)F)cc1CCl CC(C)C(=O)Nc1ccc(F)c(C(F)(F)F)c1 CN(C)C(=O)Nc1c(F)c(F)cc(F)c1CF OCc1nc(C(F)(F)F)nc2c(F)cccc12	CN(C)C(=O)Nc1cc(F)cc(C(F)(F)F)c1 CC(C)(C)OC(c1cc(F)cc(F)c1)C(F)F CCOC(=N)c1c(F)cccc1CC(F)(F)F Nc1cc(C(F)(F)F)ccc1OCCC(F)C Fc1cc(OC(F)(F)F)cc(C2CCNC2)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 250, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 250, 4.0, 2.0]
CN(C)C(=O)OC1(C(F)(F)F)CCC(F)C1 CCOC(=O)c1cc(C(F)(F)F)nc(F)c1C CC(=O)NCC(O)c1c(F)c(F)cc(F)c1F OCc1c(OC(F)(F)F)ccc(F)c1C1CC1 COc1ccc(OCC(F)(F)C(F)F)c(C)c1	Oc1cc(OC(F)(F)F)cc(F)c1CCl OC(OCC(F)(F)C(F)F)c1ccsc1 COC(=O)CC(CC(F)(F)F)c1cccc1 Cc1cc(OCC(F)(F)C(F)F)CO)ccn1 COC(=O)Cc1c(F)cnc(C(F)F)c1CF
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 250, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 250, 6.0, 1.0]
Oc1nc(F)cc(CC(F)(F)C(F)(F)F)n1 FC(C(F)(F)F)C(F)(F)COC1CCCC1 FC(F)(F)C(F)(F)COC1cccc1F OC(CCC(F)(F)F)C(F)(F)C1CC1 Cn1nc(C(F)(F)F)c(C(F)F)c1CO	CC(OC(F)(F)F)c1cccc1C(F)(F)F Cc1ncc(OC(F)(F)F)nc1C(F)(F)F OC(CCC(F)(F)C(F)(F)C(F)F)C1CC1 Fc1ccc(OCC(F)(F)F)cc1C(F)F NC(c1ccoc1)C(C(F)(F)F)C(F)(F)F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 250, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 250, 6.0, 2.0]
COC(=O)CCCC(F)(F)C(F)(F)C(F)F OC(O)c1cc(C(F)(F)F)c(F)c(F)c1F Oc1cc(C(F)(F)F)cc(C(F)(F)F)c1O COC(=O)CCCC(F)(F)C(F)(F)C(F)(F)F -- invalid--	OC(O)(CCCC(F)(F)F)CC(F)(F)F CC(=O)NCC(O)(C(F)(F)F)C(F)(F)F CCOc1ccc(C(F)(F)F)c(C(F)(F)F)c1 OC(OCC(F)(F)F)c1cccc(F)c1F Oc1c(OC(F)(F)F)cccc1C(F)(F)F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 300, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 300, 4.0, 1.0]
Cc1noc(-c2ccc(C(F)(F)F)cc2Cl)c1C(F)F Cc1cc(OC(F)(F)F)cc(F)c1 COc1cnc(C(F)(F)F)c(F)c1CBr Nc1c(F)cc(Oc2ccc(C(F)(F)F)cc2)c(Cl)c1 COc1cnc(-c2ccc(C(F)(F)F)cc2)c(Cl)c1	Fc1cnc(OC(F)(F)F)c(l)c1 O=C(Nc1ncc(C(F)(F)F)cc1Cl)c1cccc1 Cn1cc(-c2noc(-c3cc(F)c(F)c(F)c3F)n2)s1 NC(=O)c1c(F)cccc1Nc1cc(C(F)(F)F)ccn1 OC(Cc1ccc(C(F)(F)F)c(F)c1)c1cccs1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 300, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 300, 4.0, 2.0]
OC(Cc1cc(F)c(Br)cc1F)C(F)(F)O COC(=O)CCc1ccc(Cl)cc1C(F)(F)C(F)F CCOC(=O)Nc1c(C(F)(F)F)cc(F)nc1CCl N[C@@H](Cc1cc(F)c(F)c(F)c1F)c1ccc(O)cc1O CCC(=O)NCC(=O)Nc1cc(C(F)(F)F)cc(F)c1C	NC(=O)COc1cccc1-c1c(F)c(F)cc(F)c1F OC(O)(c1cc(F)c(F)c(F)c1)c1ccc(F)cc1Cl NCc1cnc(OC(F)(F)F)nc1Oc1c(F)cccc1 COc1ccc(CNc2cc(F)c(F)c(F)c2F)cc1O FC(F)(F)Oc1cccc(Oc2cc(F)cc(Cl)c2)c1



['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 300, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 300, 6.0, 1.0]
COC(c1cc(C(F)(F)F)nc(C(F)(F)F)c1)C1CC1 Nc1ccc(C(=O)NCC(F)(F)C(F)(F)C(F)F)cc1 NCc1cc(OC(F)(F)F)c(Cl)cc1C(F)(F)F NC(=O)c1cc(C(F)(F)F)cc(C(F)(F)F)c1CC1 CC(C)C(=O)Nc1cc(C(F)(F)F)cc(C(F)(F)F)c1	COc1c(C(F)(F)F)nc(C(F)(F)F)c1CC1 OC(c1ccc(F)cc1)c1c(F)c(F)c(F)c(F)c1F Nc1cc(C(F)(F)F)cc(OC(F)(F)F)c1CC1 FC(F)(F)c1cccc(-c2ccc(OC(F)(F)F)cc2)c1 Nc1c(F)cccc1Oc1cc(F)c(F)c(C(F)(F)F)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, 0.0, 300, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, 0.0, 300, 6.0, 2.0]
OCc1ccc(C(F)(F)F)cc1OCCC(F)(F)CF OB(O)c1c(C(F)(F)F)ccc(Cl)c1C(F)(F)F COc1ccc(OC(C(F)(F)F)C(F)(F)C(F)F)nc1 OCc1c(OC(F)(F)F)nc(C(F)(F)F)c1C1CC1 CS(=O)(=O)Nc1cc(C(F)(F)F)cc(C(F)(F)F)c1	COc1cc(C(F)(F)F)c(C(F)(F)F)c(CC(N)=O)c1 Cc1cc(OC(F)(F)F)nc(OC(F)(F)F)c1CC#N Cc1ccc(COCC(F)(F)C(F)(F)C(F)F)cc1O O=C(O)Cc1cc(C(F)(F)F)nc(C(F)(F)F)c1CN OC(c1cccc(OC(F)(F)F)c1)c1ccc(F)cc1F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 250, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 250, 4.0, 1.0]
Oc1ccc(-c2ccc(F)c(F)c2)c(F)c1F COC(c1c(F)c(F)nc(F)c1F)C1CCC1 NCc1cn(CC(F)(F)F)nc1CC(=O)F NCc1cc(OC(F)(F)F)cc(Cl)c1F Cn1cnc(OC(F)(F)F)c1-c1ccc(F)cc1	Cc1nc(-c2cccc(C(F)(F)F)c2F)c(C)o1 CC(NC(=O)C(F)(F)C(F)F)c1cccc1 CC(CO)Nc1cccc1C(F)(F)C(F)F CC(C)Oc1nc(C(F)(F)F)c(F)cc1CN OCC1Cc2cc(C(F)(F)F)cc(F)c2S1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 250, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 250, 4.0, 2.0]
COc1cc(C(F)(F)F)nc(OC)c1CCF CC(O)c1c(OC(F)(F)F)cc(F)cc1CN OCCC(=O)Nc1ccc(F)c(C(F)(F)F)c1 CC(N)(C(=O)O)c1nc(C(F)(F)F)ccc1F O=C(O)c1ccn(CCC(F)(F)C(F)F)c1F	OCC(O)Cc1ncc(C(F)(F)F)cc1CF O=C(O)CC(CC(F)(F)C(F)F)c1ccc[nH]1 CCC(NCC(F)(F)C(F)F)C(=O)OCC COc1ccc2c(F)c(F)c(F)c(F)c2c1O NCc1cc(OC(F)(F)F)c(O)c(CF)c1C
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 250, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 250, 6.0, 1.0]
OCc1nc(C(F)(F)F)c(C(F)(F)F)s1 OCCc1c(F)c(F)c(C(F)(F)F)c(F)c1 C[C@@H](O)c1c(F)c(F)c(C(F)(F)F)c(F)c1 OC(c1cc(F)cc(F)c1)C(F)(F)C(F)F CCc1c(OC(F)(F)F)n[nH]c1C(F)(F)F	N[C@@H](CO)c1c(F)c(F)c(F)c1CF CCc1ccc(OC(F)(F)F)c(C(F)(F)F)c1 Nc1cnc(OC(F)(F)F)c(C(F)(F)F)c1 OC[C@@H](c1cc(F)c(F)c(F)c1)C(F)(F)F COc1ncc(C(F)(F)F)c(C(F)(F)F)n1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 250, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 250, 6.0, 2.0]
C=C(O)CC(=O)C(C(F)(F)F)C(F)(F)CF OC(F)(F)C(F)(F)Oc1ccc(F)c(F)c1 OC[C@@H](O)CCC(F)(F)C(F)(F)C(F)(F)F OCCOCCC(F)(F)C(F)(F)C(F)(F)F OCc1c(O)c(F)c(F)c(F)c1C(F)(F)F	C[Si](C)(O)OC(F)(F)C(F)(F)C(F)(F)F COCC(O)CN(C(F)(F)F)C(F)(F)F CCOC(C)(O)(C(F)(F)F)C(F)(F)CF O=C(O)CCCC(F)(F)C(F)CC(F)(F)F --invalid--

['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 300, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 300, 4.0, 1.0]
O=C(Cn1nc(C(F)(F)F)cc1Cl)c1cccc1F CC1CCN(C(=O)Nc2cc(F)cc(C(F)(F)F)c2)CC1 OC(c1cc(F)cc(F)c1)c1ccc(C(F)(F)Cl)cc1 OC(c1cc(F)cc(F)c1)c1cnc(C(F)(F)Cl)cc1 CCc1nc(-c2ccc(OC(F)(F)F)cc2)nc(C)c1F	Cc1ccc(CC(=O)Nc2cc(F)cc(F)c2)c(F)c1F Cc1cc(C(F)(F)F)nc(Oc2cc(F)cc(Cl)c2)n1 CCCc1ncc(C(F)(F)F)c(Oc2ccc(F)cc2)n1 CC(OCC(F)(F)C(F)F)c1ccc(Cl)cc1Cl N#Cc1ccc(OCC(F)(F)C(F)F)cc1Br
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 300, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 300, 4.0, 2.0]
Cc1cc(OC(F)(F)C(F)(F)C(=O)NC2CC2)cs1 COC(=O)c1ncc(C(F)(F)F)c(F)c1Br COC(=O)Cc1nc(C(F)(F)F)c(F)cc1CCI CCOC(=O)Cc1cc(C(F)(F)F)cc(F)c1CCI FC(F)(F)Oc1cc(OC(F)F)cc(Br)n1	NC(=O)COc1ccc(C(F)(F)F)c(F)c1Br COC(=O)c1ccc(C(F)(F)F)c(-c2ccc(F)cc2)c1 FC(F)(F)Oc1ccc(OCC2ccncc2)c(F)c1C CCOc1cc(OC(F)(F)F)c(F)cc1Br N[C@@H](CC(=O)O)c1c(C(F)(F)F)cc(F)cc1CCI
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 300, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 300, 6.0, 1.0]
FC(F)C(F)(F)Oc1cc(C(F)(F)F)cnc1CCI OCc1c(C(F)(F)F)ccc(C(F)(F)F)c1CCI CCc1nc(OC(F)(F)F)c(C(F)(F)F)cc1CC#N OC(c1nc2cccc2s1)C(F)(F)C(F)(F)CF Fc1ccc(-c2ccc(OC(F)(F)F)cc2)c(F)c1F	Nc1ccc(OCC(F)(F)C(F)(F)C(F)F)cc1C#N Oc1cc(F)c(-c2ccc(C(F)(F)F)cc2)c(F)c1F CCN(CC(F)(F)C(F)(F)C(F)(F)F)C(=O)NC1CC1 Fc1ccc(C(F)(F)F)c(Oc2cccc(F)c2)c1F Fc1cc(OC(F)(F)F)ccc1-c1ccc(F)c(F)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.0, -0.1, 300, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.0, -0.1, 300, 6.0, 2.0]
Cc1c(OC(F)(F)F)cnc(C(F)(F)F)c1CC(N)=O O=C(O)C(CC(F)(F)C(F)(F)C(F)F)c1cccc1 OCc1cc(OC(F)(F)F)c(Cl)cc1C(F)(F)F C[C@@](N)(C(=O)O)c1nc(C(F)(F)F)c(C(F)(F)F)n1C CCOC(=O)c1cc(C(F)(F)F)ccc1C(F)(F)CF	Oc1ccc(COCC(F)(F)C(F)(F)C(F)F)cc1F OC(O)(Cc1ccc(C(F)(F)F)cc1)CC(F)(F)CF OC(O)(c1cc(C(F)(F)F)cc(C(F)(F)F)c1)C1CC1 O=C(NCC(F)(F)C(F)(F)C(F)F)c1ccc(O)cc1 FC(F)(F)COCCOc1c(F)cc(C(F)F)cc1N
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 250, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 250, 4.0, 1.0]
Cc1ccc(CNC(=O)C(F)(F)C(F)F)cc1 Cc1cc(CC(=O)NCC(F)(F)F)ccc1F O=C(c1cccnc1)c1cc(F)c(F)c(F)c1F FC(F)(F)c1cccc(Oc2cccc2)c1F CCNC(=O)Nc1ccc(F)c(C(F)(F)F)c1	Cc1c(OC(F)(F)F)cc(F)cc1CCI OC(CC(F)(F)F)c1ccc(F)c(Cl)c1 COc1cc(C(F)(F)F)c(F)cc1CCI OCc1cnc(C(F)F)c(Cl)c1C(F)(F)F Nc1cc(OCC(F)(F)C(F)F)ccc1C#N
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 250, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 250, 4.0, 2.0]
OCCC(=O)Nc1c(F)cccc1C(F)(F)F O=C(Cc1ccc(OC(F)(F)F)cc1)C1CC1 CCOc1c(OC(F)(F)F)cc(F)cc1C#N CCOC1(C(F)(F)F)Oc2ccc(F)cc2C1 CC(C)CC(=O)NCC(O)CC(F)(F)C(F)F	OCc1cc(OCCC(F)(F)F)cc(F)c1C OCc1c(C(F)F)ncc(OC(F)F)c1CC COc1nc(OC(F)(F)F)c(F)cc1CC#N Cc1ccc(C(F)(F)C(F)(F)C(=O)O)cc1C Cc1nc(C(F)(F)F)c(CC(=O)O)cc1CF

['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 250, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 250, 6.0, 1.0]
O=C(NCC(F)(F)C(F)CF)CC(F)(F)F Cc1cnc(OC(F)(F)F)c(C(F)(F)F)c1 O=C(Nc1cc(F)cc(F)c1F)C(F)(F)F NCC(O)CCC(C(F)(F)F)C(F)(F)F Fc1cc(OCC(F)(F)F)ccc1C(F)F	OC(c1cc(F)c(F)c(F)c1)CC(F)(F)F FCOc1cc(C(F)(F)F)ccc1C(F)(F)F Cc1ncc(C(F)(F)F)c(OC(F)(F)F)n1 OC(c1c(F)c(F)c(F)c(F)c1F)C(F)F OCCc1c(F)c(F)c(C(F)(F)F)c(F)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 250, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 250, 6.0, 2.0]
O=C(O)CCCC(F)(F)C(F)(F)C(F)F CCOC(=O)CC(C(F)(F)F)C(F)(F)CF Fc1cccc(OC(F)(F)OC(F)(F)F)c1 Oc1cc(OC(F)(F)F)cc(C(F)(F)F)c1 CCC(=O)OCC(C(F)(F)F)C(F)(F)CF	FC(F)(F)C1(C(F)(F)F)OCCCCOC1 COC(=O)CCC(C(F)(F)F)C(F)(F)CF Cc1c(O)cc(C(F)(F)F)cc1OC(F)(F)F O=C(OCCCC(F)(F)F)CC(F)(F)F CC(O)CCC(O)(C(F)(F)F)C(F)(F)F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 300, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 300, 4.0, 1.0]
O=C(NCC(F)(F)C(F)F)c1ccc(Cl)cc1Cl CCc1ncc(OC(F)(F)F)c(F)c1CBr FC(F)(F)c1ccc(Oc2ncccc2Cl)c(F)c1 Nc1ncc(F)cc1C(=O)Nc1ccc(C(F)(F)F)cc1 CCCc1nc(OC(F)(F)F)nc(F)c1Br	CN(C)c1cc(C(F)(F)F)nc(Oc2ccc(F)cc2)n1 OC(c1cccc(C(F)(F)F)c1)c1c(F)cccc1Cl Cc1cnc(C(F)(F)F)c(Oc2ccc(F)c(Cl)c2)n1 N#Cc1ccc(COc2cccc2C(F)(F)F)c(F)c1 CNc1ccc(Oc2ccc(C(F)(F)F)nc2)c(F)c1C
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 300, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 300, 4.0, 2.0]
CC(Oc1cccc(C(F)(F)F)c1)C(=O)c1cccc1 COc1nc(OC(F)(F)F)c(F)cc1Br O=C(NCC(F)(F)CO)c1c(F)cc(F)cc1CCl FC(F)(F)Oc1cc(F)c(OCc2cccnc2)c(C)c1 CS(=O)(=O)Nc1ccc(C(F)(F)C(F)F)c(Cl)c1	CS(=O)(=O)N1CCN(c2c(F)c(F)cc(F)c2F)C1 Oc1cc(OC(F)(F)F)ccc1CBr CS(=O)(=O)Nc1ccc(SC(F)(F)C(F)F)cc1 O=S(=O)(c1cccc(C(F)(F)F)c1)c1ccc(F)cc1 CC(Nc1cc(C(F)(F)F)cc(F)c1)C(=O)OC(C)(C)C
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 300, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 300, 6.0, 1.0]
Cc1c(C(F)(F)F)cc(OC(F)(F)F)nc1CCl Nc1cc(OC(F)(F)F)cc(C(F)(F)F)c1CCl FC(F)(F)Oc1ccc(-c2ccc(C(F)(F)F)cc2)cc1 O=Cc1ccc(C(F)(F)F)c(C(F)(F)F)c1CCl --invalid--	Oc1c(F)cccc1-c1cc(C(F)(F)F)cc(F)c1F OC(c1cccc(C(F)(F)F)c1)c1cccc(F)c1F CC(=O)Nc1c(C(F)(F)F)cnc(C(F)(F)F)c1CN OCc1cc(C(F)(F)F)cc(C(F)(F)F)c1CCl OC(c1ccc(C(F)(F)F)nc1)c1c(F)cc(F)cc1F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, 0.0, 300, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, 0.0, 300, 6.0, 2.0]
OCC(Oc1cccc(C(F)(F)F)c1)CC(F)(F)CF O=Cc1cc(OCC(F)(F)C(F)(F)C(F)F)ccc1N OCc1c(OC(F)(F)F)ncc(C(F)(F)F)c1Cl NC(COCC(F)(F)F)c1cccc(OC(F)(F)F)c1 OCc1cccc1OCCC(F)(F)C(F)(F)C(F)F	C[C@H](NC(=O)O)c1cc(C(F)(F)F)cc(C(F)(F)F)c1 CCOC(=O)Cc1c(F)cc(C(F)(F)F)cc1C(F)F NC(=O)c1c(OC(F)(F)F)cnc(C(F)(F)F)c1CN O=C(NCC(F)(F)F)N1CCC(O)(C(F)(F)F)CC1 FC(F)(F)Oc1ccc(OC(F)(F)C(F)Cl)cc1

['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 250, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 250, 4.0, 1.0]
CNCC(=O)Nc1cc(F)c(C(F)(F)F)cc1 OC(CC(F)(F)F)c1ccc(F)c(Cl)c1 Fc1ccc(COC2CCCC2F)c(F)c1F OCc1cc(C(F)(F)F)c(F)cc1CCI COC1c(C(F)(F)F)ccc(F)c1CCI	OC(c1cc(F)cc(C(F)(F)F)c1)C1CCC1 OCCc1cc(C(F)(F)F)c(Cl)cc1F OCCc1nc(C(F)(F)F)c(F)cc1Cl Nc1cnc(OCCCC(F)(F)F)c(F)c1F Nc1ccc(OCCC(F)(F)C(F)F)c(C)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 250, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 250, 4.0, 2.0]
CN(C(=O)O)c1cc(C(F)(F)F)cc(F)c1C CS(=O)(=O)CCSCC(F)(F)C(F)F CCC(=O)NCC(=O)NCC(F)(F)C(F)F CNC1CC(C(F)(F)C(F)(F)C(=O)O)CC1 COC(=O)Nc1cccc(C(F)(F)C(F)F)c1	CCCC(=O)OCCCC(F)(F)C(F)(F)F O=Cc1cc(OCC(F)(F)C(F)F)cs1 OC(COCC(F)(F)F)c1ccc(F)c(C)c1 CCCOC(=O)Nc1c(F)c(F)cc(F)c1F CCOC(=O)Nc1nc(C(F)(F)F)ccc1F
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 250, 6.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 250, 6.0, 1.0]
C[C@H](O)c1cc(C(F)(F)F)cc(C(F)(F)F)c1 OCC(F)(F)C(F)(F)c1cc(F)cc(F)c1 C[C@H](N)CC(=O)NC(C(F)(F)F)C(F)(F)F Cc1cc(OC(F)(F)F)nc(C(F)(F)F)c1 --invalid--	Cc1c(OC(F)(F)F)cccc1C(F)(F)F Nc1ncc(F)c(OC(F)(F)F)c1C(F)F CNCC(=O)NCC(F)(F)C(F)(F)C(F)F Fc1cc(F)c(OCCC(F)(F)F)c(F)c1 FC(F)(F)CCOc1cc(F)c(F)c(F)c1
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 250, 6.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 250, 6.0, 2.0]
OCc1c(F)c(F)c(OC(F)(F)F)c(F)c1 C=CCOC(=O)CC(F)(F)C(F)(F)C(F)F OCc1cc(C(F)(F)F)cc(C(F)(F)F)c1O Cc1c(OC(F)(F)F)[nH]c(C(F)(F)F)c1O C=CCOC(=O)C(C(F)(F)F)C(F)(F)CF	COC(=O)CCC(C(F)(F)F)C(F)(F)CF CC(CC(F)(F)C(F)(F)C(F)F)CC(=O)O CC(C)(CC(F)(F)C(F)(F)C(F)F)C(=O)O Cc1ccc(OC(F)(F)F)c(OC(F)(F)F)c1 FCOc1c(C(F)(F)F)[nH]c(C(F)F)c1O
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 300, 4.0, 1.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 300, 4.0, 1.0]
OC(c1c(F)c(F)c(F)c1Br)C1CC1 Oc1nc(F)c(C(F)(F)F)cc1 OC(Cc1ccc(F)c(Br)c1)CC(F)(F)F COC1nc(C(F)(F)F)ccc1CBr O=C(Cc1nc(C(F)(F)F)ns1)c1ccc(F)cc1C	OC(c1cccc(C(F)(F)F)c1)c1ccc(F)cc1Cl CC1CC(c2ccc(C(F)(F)F)cc2)NC(=O)NC1F CC(NC(=O)CC(F)(F)C(F)F)c1ccc(Cl)cc1 COC1cc(Br)cc(CC(F)(F)C(F)F)c1 Cn1nc(C(F)(F)F)cc1C(=O)Nc1ccc(F)cc1C
['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.5, 7.5, -0.1, 300, 4.0, 2.0]	['EA', 'IE', 'LogVis', 'MolWt', 'n_F', 'n_O'] : [0.0, 7.5, -0.1, 300, 4.0, 2.0]
CC(=O)Nc1nc(-c2ccc(C(F)(F)F)c(F)c2)c(C)o1 Oc1cccc(Oc2c(F)c(F)c(F)c2Cl)c1C O=S(=O)(Cc1ccc(F)cc1)c1c(F)cc(F)cc1F OC(O)(c1c(F)cccc1F)c1cc(F)c(Cl)cc1F --invalid--	CC(Oc1cccc(C(F)(F)F)c1)c1cccc(F)c1O COC1nc(C(F)(F)F)ccc1OCc1ccc(F)cc1 OC(COCC(F)(F)C(F)F)Cc1ncccc1Cl FC(F)C(F)(F)Oc1cccc1COC1cccc1 COC(=O)c1ncc(C(F)(F)F)cc1-c1ccc(F)cc1

