

## 1 **AGILE Platform: A Deep Learning-Powered Approach to Accelerate LNP Development** 2 **for mRNA Delivery**

3 Yue Xu<sup>1#</sup>, Shihao Ma<sup>4,5,6#</sup>, Haotian Cui<sup>4,5,6#</sup>, Jingan Chen<sup>2</sup>, Shufen Xu<sup>1</sup>, Kevin Wang<sup>1</sup>, Andrew  
4 Varley<sup>1</sup>, Rick Xing Ze Lu<sup>2</sup>, Bo Wang<sup>4,5,6,7\*</sup> and Bowen Li<sup>1,2,3\*</sup>

5 <sup>1</sup> Department of Pharmaceutical Sciences, Leslie Dan Faculty of Pharmacy, University of  
6 Toronto, Toronto, Ontario, Canada.

7 <sup>2</sup> Institute of Biomedical Engineering, University of Toronto, Toronto, Ontario, Canada.

8 <sup>3</sup> Princess Margaret Cancer Center, University Health Network, Toronto, ON M5G 2C1, Canada

9 <sup>4</sup> Department of Computer Science, University of Toronto, Toronto, ON, Canada.

10 <sup>5</sup> Vector Institute for Artificial Intelligence, Toronto, ON, Canada.

11 <sup>6</sup> Peter Munk Cardiac Centre, University Health Network, Toronto, ON, Canada.

12 <sup>7</sup> Department of Laboratory Medicine and Pathobiology, University of Toronto, Toronto, ON,  
13 Canada.

14 # These authors contributed equally.

15 \* Communication can be sent to [bowang@vectorinstitute.ai](mailto:bowang@vectorinstitute.ai); [bw.li@utoronto.ca](mailto:bw.li@utoronto.ca)

### 16 **Abstract**

17 Ionizable lipid nanoparticles (LNPs) have seen widespread use in mRNA delivery for clinical  
18 applications, notably in SARS-CoV-2 mRNA vaccines. Despite their successful use, expansion of  
19 mRNA therapies beyond COVID-19 is impeded by the absence of LNPs tailored to different target  
20 cell types. The traditional process of LNP development remains labor-intensive and cost-  
21 inefficient, relying heavily on trial and error. In this study, we present the **AI-Guided Ionizable**  
22 **Lipid Engineering (AGILE)** platform, a synergistic combination of deep learning and  
23 combinatorial chemistry. AGILE streamlines the iterative development of ionizable lipids, crucial  
24 components for LNP-mediated mRNA delivery. This approach brings forth three significant  
25 features: efficient design and synthesis of combinatorial lipid libraries, comprehensive in silico  
26 lipid screening employing deep neural networks, and adaptability to diverse cell lines. Using  
27 AGILE, we were able to rapidly design, synthesize, and evaluate new ionizable lipids for mRNA  
28 delivery in muscle and immune cells, selecting from a library of over 10,000 candidates.  
29 Importantly, AGILE has revealed cell-specific preferences for ionizable lipids, indicating the need  
30 for different tail lengths and head groups for optimal delivery to varying cell types. These results  
31 underscore the potential of AGILE in expediting the development of customized LNPs. This could  
32 significantly contribute to addressing the complex needs of mRNA delivery in clinical practice,  
33 thereby broadening the scope and efficacy of mRNA therapies.

### 34 **One Sentence Summary**

35 AI and combinatorial chemistry expedite ionizable lipid creation for mRNA delivery.

36

37 .

38

39

## 40 Introduction

41 Messenger RNA (mRNA) has emerged as a versatile tool with wide-ranging biomedical  
42 applications, ranging from vaccines and protein replacement therapy to cell engineering and gene  
43 editing<sup>1,2</sup>. This versatility has fueled widespread interest in exploiting mRNA to tackle an array  
44 of diseases<sup>3,4</sup>. However, the inherently unstable nature of mRNA and its susceptibility to nuclease  
45 degradation necessitates an effective delivery system, a role typically fulfilled by ionizable lipid  
46 nanoparticles (LNPs)<sup>5</sup>. Both Comirnaty and Spikevax, two SARS-CoV-2 vaccines approved  
47 amidst the COVID-19 pandemic, are grounded on LNP-based mRNA delivery<sup>6,7</sup>. Moreover, LNP  
48 technology helped the first siRNA drug (Onpatro) obtain U.S. FDA approval in 2018<sup>8-10</sup>. The  
49 classical LNP formulation comprises four compositions: ionizable lipids, cholesterol, helper lipids,  
50 and PEGylated lipids. Notably, each of the three FDA-approved RNA LNPs has a distinct  
51 ionizable lipid design, highlighting the pivotal role of ionizable lipids in LNP technology. Their  
52 primary functions include packaging mRNA into LNPs and facilitating its entry into the cytoplasm  
53 of target cells for ribosomal binding and subsequent protein expression<sup>11-14</sup>. An ionizable lipid  
54 generally consists of an ionizable amine head group and two lipid tails. This structure enables  
55 protonation at acidic pH, thereby adopting a cationic character during the LNP formulation process,  
56 facilitating the encapsulation of anionic RNA molecules. At physiological pH, the ionizable lipid  
57 remains neutrally charged, thereby circumventing potential toxicity associated with non-ionizable  
58 cationic lipids. Once the LNP encapsulating mRNA is endocytosed, ionizable lipids undergo  
59 protonation again in the acidic endosomal environment, disrupting the inner phospholipid  
60 membrane of endosomes and promoting the release of mRNA into the cytoplasm of target cells.  
61 As the COVID-19 pandemic recedes, the spectrum of mRNA applications continues to broaden  
62 beyond vaccination, thus emphasizing the necessity for a diverse array of ionizable lipids  
63 proficient in mRNA delivery to a variety of target cells and tissues.

64 Although previous research has provided some insight into the rational design of ionizable lipids  
65 to improve the mRNA delivery performance of LNPs, the approach often covers limited structural  
66 space, potentially overlooking some promising lipid designs. Combinatorial chemistry, employing  
67 multi-component reactions, has recently been used to enable high-throughput synthesis (HTS) of  
68 extensive and chemically diverse lipid libraries. For example, a Ugi-based three-component  
69 reaction (3-CR) could enable the swift synthesis of a combinatorial library, comprising 1,080  
70 ionizable lipids, ultimately leading to the identification of a STING-activating ionizable lipid  
71 conducive to mRNA vaccine delivery<sup>15</sup>. More recently, another 3-CR system based on the  
72 Michael addition was used to generate a library of over 700 ionizable lipids, resulting in the  
73 discovery of a potent lipid uniquely suited for efficient mRNA delivery to the lung epithelium<sup>16</sup>.  
74 While the 3-CR combinatorial chemistry has been showcased to facilitate the synthesis of new  
75 ionizable lipids, constructing and testing a more extensive lipid library, running into hundreds of  
76 thousands of compounds, for mRNA transfection in different cell targets remains a formidable,  
77 time-consuming, and costly task<sup>17</sup>. This challenge consequently restricts efforts to design and test  
78 more diverse and innovative structures. New strategies are essential to hasten the discovery and  
79 optimization of ionizable lipids for achieving desirable mRNA transfection in specific target cells.

80 Deep learning, a subset of artificial intelligence (AI), poses a promising resolution to the challenge  
81 of exploring molecular search spaces<sup>18-20</sup>. With ample high-quality training data, these techniques

82 can effectively extract insights from observed molecules, capitalizing on underlying chemical  
83 structures and properties, and extrapolating to a broader array of unobserved molecules. Indeed,  
84 the rise of deep learning is reshaping chemical compound discovery, transforming this process  
85 from a trial-and-error practice to an intelligent, data-driven strategy <sup>21-27</sup>. In this study, we  
86 pioneered utilizing cutting-edge deep learning methodologies to accelerate the development of  
87 ionizable lipids for mRNA delivery, culminating in the **AI-Guided Ionizable Lipid Engineering**  
88 (AGILE) platform. This platform not only dramatically expands the molecular space of lipid  
89 structures by several magnitudes, but also significantly truncates the timeline for new ionizable  
90 lipid development, reducing it from potential months or even years to weeks. Essentially, AGILE  
91 employs a pre-trained deep-learning neural network that assimilates structural knowledge from  
92 millions of small-molecule components. The model utilizes vast amounts of unlabeled data from  
93 a combinatorial lipid library, employing a self-supervised approach to learn differentiable lipid  
94 representations. Following the fine-tuning on wet-lab data collected after HTS, AGILE can  
95 identify promising lipids for high mRNA transfection potency in specific cells from a significantly  
96 larger combinatorial library with enhanced accuracy. Leveraging this workflow, we fine-tuned the  
97 deep learning model using transfection data from Hela cells, which subsequently led to the  
98 prediction of 15 top lipid structures from a pool of 12,000 lipid candidates. This process facilitates  
99 the identification of an ionizable lipid H9 that shows superior mRNA transfection potency  
100 compared to LNPs containing (D-Lin-MC3-DMA) <sup>2</sup>, an FDA-approved ionizable lipid for RNA  
101 delivery, following intramuscular injection. Notably, the transfection effect of H9 LNPs is  
102 localized to the muscle, with significantly less off-target transfection in other tissues, such as the  
103 liver. Moreover, we showed that AGILE could be quickly repurposed to discover LNPs for other  
104 target cells, as demonstrated by identifying a new lipid, R6, optimized for mRNA delivery to  
105 macrophages. Experimental observations, such as the significance of non-biodegradable tail  
106 structures in macrophage transfection and the correlation between the carbon chain length and  
107 transfection potency, underscore AGILE's potential to provide meaningful biological insights and  
108 tailor LNPs for individual cell types. AGILE's ability to customize for different cell types suggests  
109 its potential to steer the formulation of new mRNA-LNPs, finely tailored to various clinical  
110 scenarios.

## 111 Results

### 112 Overview of the AGILE platform.

113 By synergistically integrating deep learning methodologies with combinatorial lipid synthesis  
114 chemistry, AGILE is dedicated to streamlining the discovery process for new ionizable lipids,  
115 which are crucial to LNP-based mRNA delivery. Central to this platform is a suite of deep learning  
116 algorithms, collectively referred to as the AGILE model. This model, encompassing a graph  
117 encoder and a molecular descriptor encoder, adeptly captures the intrinsic characteristics of  
118 ionizable lipid molecular structures and their corresponding chemical attributes. The  
119 implementation of AGILE in this study unfolds over three key stages, as illustrated in Figure 1a:  
120 (1) the constitution of a virtual library and initial self-supervised model training, (2) the acquisition  
121 of empirical data from an experimental library, enhancing the precision of the pre-trained model  
122 through supervised fine-tuning, and (3) the execution of *in silico* analysis on ionizable lipids in a  
123 candidate library, leveraging the refined deep learning algorithms (Methods 1.1 for additional  
124 details). As a multifaceted tool, AGILE generates predictions on the mRNA transfection capacity  
125 of ionizable lipids in LNP formulations and significantly facilitates the design of LNP for specific  
126 target cells.

127 Stage 1 aims to develop a graph encoder proficient in differentiating and depicting distinct lipids  
128 through pre-training on a vast collection of unlabeled lipid molecules (Methods 1.3 and 1.4). This  
129 process begins with the construction of a graph encoder utilizing Graph Neural Networks (GNN),  
130 primed with parameters from the MolCLR model, which has undergone pre-training on a  
131 repertoire of over 10 million small molecules. This “warm-starting” strategy, embedding general  
132 knowledge of small molecular structures into our algorithm, fortifies the accuracy of AGILE in  
133 subsequent stages (Supplementary Fig. S1). The graph encoder subsequently underwent  
134 continuous pre-training on a virtual library of 60,000 chemically diverse lipids through contrastive  
135 learning<sup>28</sup>, enabling the differentiation of atoms and bonds in each molecule, and thus capturing  
136 the disparities amongst various lipid structures (see Methods 1.4). This virtual library, composed  
137 of lipids with diverse amine head groups and two unique alkyl chains (Fig. 1b), is designed based  
138 on 3-CR chemistry principles, thus amenable to high-throughput combinatorial synthesis<sup>29</sup>.  
139 Overall, the pre-training in Stage 1 equips the graph encoder with a comprehensive understanding  
140 of lipid structures, thereby enhancing subsequent steps (Supplementary Fig. S1). Stage 2 seeks to  
141 further train the AGILE model with mRNA transfection potency data from a pool of ionizable  
142 lipids. To this end, we synthesized 1200 ionizable lipids by 3-CR and assessed their transfection  
143 potency in a target cell line, from which the data was leveraged to fine-tune the AGILE model in  
144 a supervised manner (Methods 1.5). To enhance the generalizability and precision, we added a  
145 molecular descriptor encoder that takes molecular descriptors computed by Mordred as the input  
146<sup>30</sup> (Methods 1.3). The output of the molecular descriptor encoder was utilized to update the  
147 representation of lipid structures by the pre-trained graph encoder. As such, the AGILE model has  
148 been trained to minimize the difference between the predicted result and the ground truth from  
149 wet-lab experiments during the fine-tuning process. Prior to the *in silico* screening in Stage 3, we  
150 assembled a candidate library containing 12,000 lipid structures by rationally selecting structures  
151 from the virtual library in Stage 1 (Fig. 1c) following three rules (Methods 1.1): (1) Removal of  
152 non-ionizable cationic lipids due to the potential risk of toxicity<sup>31</sup>; (2) Removal of lipids with too

153 short (<C10) or too long (> C18) alkyl chains based on empirical experience<sup>15</sup>; and (3) Removal  
154 of lipids requiring unavailable reagents for synthesis. The fine-tuned AGILE model was then  
155 utilized to predict the mRNA transfection potency of lipids in the candidate library, followed by a  
156 head and tail-wise ranking methodology to increase the structural diversity of the top-ranked  
157 candidates (Fig. 1d, Methods 1.6). Based on the information afforded by AGILE, the top-ranked  
158 ionizable lipid structures were selectively synthesized in the wet lab and formulated into LNPs for  
159 validating their ability to efficiently deliver mRNA to a specific target cell.

160

### 161 **Combinatorial Lipid library synthesis and screening for fine-tuning.**

162 Upon completing the pre-training of the entire virtual library in Stage 1, we tailored the model for  
163 transfection potency prediction through supervised fine-tuning. This stage involved training the  
164 model based on *in vitro* screening results, enabling the model to capture the potential transfection  
165 ability of molecules. To rapidly generate ionizable lipid libraries with high chemical diversity, we  
166 developed an automated high-throughput synthesis (HTS) platform based on the one-pot Ugi 3CR  
167 (Fig. 2a and Supplementary Fig. S2), which enabled the synthesis of a large batch (1,200) of  
168 ionizable lipids within 24 hours. The synthesized lipid library comprises 20 diverse head groups,  
169 12 alkyl chains with biodegradable ester linkages, and 5 alkyl chains containing isocyanide  
170 function groups (Fig. 2b)<sup>32</sup>. Using the HTS platform, we formulated LNPs via a liquid handling  
171 robot following a previously established classical four-composition formulation ratio<sup>33</sup>.

172 The LNPs were subsequently synthesized for testing lead candidates by four classic formulations  
173 with the ionizable lipids, helper lipid (DOPE), cholesterol, and polyethylene glycol (PEG)-  
174 phospholipid conjugate (DMG-PEG2000) (Fig. 2c)<sup>34</sup>. To evaluate the mRNA transfection  
175 potency in HeLa cells, we measured firefly luciferase (Fluc) protein expression activity by  
176 encapsulating Fluc mRNA in the LNPs (Fig. 2d). Most of these 1200 lipids showed improved  
177 mRNA transfection potency in HeLa cells compared to untreated cells (Fig. 2e). HeLa cells are  
178 commonly utilized as an *in vitro* screening model for evaluating transfection potency through  
179 intramuscular injection. This is due to their reliable expression of the low-density lipoprotein  
180 receptor (LDLR), which plays a crucial role in the cellular uptake of lipid nanoparticles (LNPs)  
181 associated with lipoproteins in the bloodstream<sup>35</sup>. The presence of LDLR in HeLa cells allows for  
182 enhanced cellular uptake of LNPs, making them a valuable tool for assessing the effectiveness of  
183 intramuscular delivery methods. Previous studies emphasized the preference for muscle as the site  
184 of vaccination. This choice was based on the rich blood supply in muscle tissue, which enables the  
185 efficient processing of foreign antigens by immune cells, leading to a robust immune response<sup>36</sup>.  
186<sup>37</sup>. Therefore, the strong correlation between transfection potency in HeLa cells and in muscle  
187 tissues further establishes their utility in evaluating the effectiveness of intramuscular delivery  
188 methods<sup>39</sup>. Meanwhile, the potencies vary significantly among test lipids, with relative luciferase  
189 units ranging from poor ( $\text{Log}_2 < 5$ ) to outstanding performance ( $\text{Log}_2 > 10$ ) (Supplementary Fig. S3).  
190 These variations can be readily used, in the fine-tuning stage, to supervise the model to learn the  
191 relation between molecule properties and its transfection potency (Methods 1.5). We used 80% of  
192 the data for the model training, 10% for selecting the best hyperparameters, and the last 10% for

193 internal verification. We observed a constant decrease of loss value on both training and validate  
194 data (Fig. 2f) and thus used the model with the lowest validation loss.

195 To verify the quality of the predictions, we split the predicted and actual *in vitro* potency values  
196 into six equal percentiles. We visualize the precision matrix on all 1,200 lipids in Fig. 2g. Although  
197 the prediction task is extremely challenging, the model works particularly better for predicting the  
198 top and least performing lipids, which are arguably the most important and informative for  
199 selecting lipid candidates. For example, a predicted top-16% performing lipid will have a chance  
200 of 0.41 to be one of the actual top-16% performing lipids found *in vitro* (Fig. 2g). We also  
201 examined our predictions using UMAP embedding (Fig. 2h)<sup>38</sup>. The UMAP algorithm assigns  
202 close LNPs presentations to adjacent points in a two-dimensional space, which are then colored  
203 based on their predicted transfection potency. The lipids gathered into regional structures with  
204 similar potency values on the resulting UMAP plot, which verifies that the learned representations  
205 capture the potential transfection ability of lipids.

206

### 207 **AGILE predicts and identifies the efficient lipid for muscle injection.**

208 With the fine-tuned model, we perform model prediction on the candidate library to screen  
209 potential lipids for muscle injection. We visualize our predictions using UMAP (Fig. 3a), and the  
210 resulting plot shows a clear separation between high and low predicted values, indicating the  
211 robustness of the model in differentiating efficacious and less efficacious ionizable lipids in a  
212 larger screening library. A closer look at the stratified distribution plots reveals that predicted  
213 potencies are clearly sorted by head group and tail combinations (Fig. 3b and Supplementary Fig.  
214 S4). Even among the top 5 performing head groups, A8 and A21 had higher predicted potencies  
215 than the others. While the tail combinations displayed less pronounced stratification of predicted  
216 transfection potencies compared to the head groups (Fig. 3c and Supplementary Fig. S4), the top  
217 tail combinations were still essential for candidate selection compared to the bottom tail  
218 combinations. The model appeared to favor unsaturated alkyl chains, a finding that was consistent  
219 with much of the literature that had been reported (Supplementary Fig. S5)<sup>39,40</sup>. Using our ranking  
220 system, which prioritizes structural diversity among lipids by considering head groups and tail  
221 combinations (Fig. 1d), we finalized a set of 15 lipid candidates (Supplementary Fig. S5).

222 We rapidly synthesized the 15 lead candidates ranked by the model in the HTS system and  
223 evaluated them in Hela cells and found that all 15 lead candidates resulted in luciferase protein  
224 expression compared with the untreated group (Fig. 3d). To investigate their potential *in vivo*, we  
225 administered mice with Fluc mRNA encapsulated in LNPs by intramuscular injection. Among 15  
226 different candidates, we observed a notably robust bioluminescence signal for H9 LNPs  
227 (Supplementary Fig. S6). After optimizing the LNPs formulation of H9 by using the design-of-  
228 experiment (DoE) (Fig. 3e, Supplementary Table. S1 and Fig. S7). After conducting a comparison  
229 with MC3 LNPs, we discovered that H9 LNPs had 2.3 times more mRNA transfection potency  
230 than MC3, which is a benchmark ionizable lipid currently used in the clinic (Fig. 3f)<sup>41</sup>. Based on  
231 the positive outcome, we proceeded to use the H9 LNPs to assess mRNA transfection potency in  
232 mice through intramuscular injection (Fig. 3g). Our findings revealed that the transfection potency  
233 of the H9 LNPs in muscle site was 7.8 times stronger than that of the MC3, with no significant

234 difference compared to the ALC-0315 (the ionizable lipid used in the SARS-CoV-2 vaccine,  
235 BNT162b2, from BioNTech and Pfizer) (Fig. 3h and i). It is worth noting that administering  
236 mRNA LNPs through intramuscular injection may cause an off-target effect, leading to the  
237 production of FLuc protein expression in the liver of mice<sup>42</sup>. When compared to ALC-0315 LNPs,  
238 H9 LNPs were found to have lower off-target effects in the liver while maintaining similar  
239 transfection effectiveness in muscle tissue. (Fig. 3j and k). Inspired by these findings, we  
240 investigated the potential of H9 LNP for vaccination. To compare the delivering efficacy of H9  
241 and ALC0315 LNPs, we administered cre-recombinase mRNA LNPs to mTmG reporter mouse  
242 models<sup>43</sup>. These mice harbored gene mutations in the Gt(ROSA)26Sor locus, and upon cre-mRNA  
243 expression, the mT cassette was excised in the cre-expressing tissue, enabling the expression of  
244 the downstream membrane-targeted green fluorescent protein (GFP, mG) cassette (Fig. 3l). We  
245 observed comparable levels of GFP protein expression at the intramuscular injection site for H9  
246 and ALC-0315 LNPs. However, ALC-0315 LNPs showed higher protein expression levels in liver  
247 tissue (Supplementary Fig. S8). Quantification of confocal images revealed that the H9 LNP  
248 exhibits 28% lower transfection potency in the liver compared to ALC-0315 LNPs (Fig. 3m).  
249 Notably, clinical studies have associated ALC-0315-based BNT162b2 mRNA vaccines with  
250 autoimmune hepatitis (AIH) following vaccination<sup>44</sup>. Hence, it is anticipated that the H9 LNPs  
251 predicted by AGILE will alleviate the serious potential side effects of hepatitis with a lower off-  
252 targeting effect.

### 253 **Using AGILE to identify ionizable lipids for Macrophage mRNA delivery.**

254 It is known that conventional adeno-associated virus (AAV) vectors struggle to transfer innate  
255 immune cells, which highlights the importance of a non-viral mRNA delivery system<sup>45</sup>. Although  
256 non-viral delivery vectors may avoid this disadvantage in immune cells, they still require effective  
257 mRNA transfection potency into the targeted immune cell type<sup>46</sup>. In order to test AGILE's ability  
258 to identify ionizable lipids that can efficiently transfect immune cells, we examined 1,200 lipids  
259 in RAW 264.7 cells (a macrophage cell line). This allowed us to create a dataset specifically for  
260 macrophages and fine-tune the screening process. The results revealed considerable differences in  
261 transfection potency between these two cell lines, with even the same batch of lipids showing  
262 totally disparate outcomes in HeLa cells and macrophage cells (Supplementary Fig. S9). The study  
263 discovered that immune cells were less easily transfected by LNPs than HeLa cells, demonstrating  
264 that immune cells pose a greater challenge for transfection. (Supplementary Fig. S10)<sup>47, 48</sup>.

265 With the model fine-tuned on the macrophage-specific dataset, we once again performed model  
266 prediction and visualized the predicted transfection potencies for RAW 264.7 cells using UMAP  
267 (Fig. 4a). Contrasting with the UMAP of predicted potencies for HeLa cells, the top-tier predicted  
268 LNPs are dispersed more widely throughout the space, potentially suggesting an increased  
269 complexity in predicting potencies for RAW 264.7 cells. Mirroring the pattern observed in HeLa  
270 cells, the predicted potencies for RAW 264.7 cells exhibit evident stratification when categorized  
271 by head groups and tail combinations (Fig. 4b, c and Supplementary Fig. S4). The top 15  
272 candidates were synthesized in the wet lab and subjected to an initial screen in RAW 264.7 cells,  
273 where 11 out of 15 showed improved transfection potency compared to MC3 (Fig. 4d and  
274 Supplementary Fig. S11). R6 was chosen as the best-performing lipid among the 15 candidates  
275 and subjected to formulation optimization using the design of experiments (DoE) (Fig. 4e and

276 Supplementary Fig. S12)<sup>49</sup>. We then loaded LNPs with Fluc mRNA and evaluated luciferase  
277 protein expression in both RAW 264.7 and Hela cells to compare H9 and R6 LNPs performance  
278 in different cell lines. Interestingly, the results were quite different in RAW 264.7 cells, where R6  
279 exhibited significantly higher transfection potency than H9 and MC3 (Fig. 4f). However, H9  
280 demonstrated more than a 2-fold increase in transfection potency compared to R6 in Hela cells  
281 (Fig. 4g). These results demonstrated the necessity to develop LNPs specifically for individual cell  
282 types and tissues, rather than a one-size-fits-all approach for all targets. Based on the excellent  
283 performance of R6 LNPs in RAW 264.7, we tried to use R6 LNPs to deliver GFP mRNA to RAW  
284 264.7. When compared to H9 and MC3 LNPs, R6 LNPs exhibited a 5-fold increase in transfection  
285 potency in RAW 264.7 as determined by flow cytometry (Fig. 4h and 4i). These results validated  
286 the success of AGILE in identifying a new ionizable lipid for efficient macrophage transfection,  
287 highlighting its potential to be utilized for the development of non-viral mRNA delivery vectors  
288 for immune cells.

### 289 **Interpretation of the AGILE deep learning model.**

290 AGILE elucidates its models through two mechanisms: (1) identification of influential molecular  
291 descriptors using a gradient-based model interpretation method, and (2) discernment of critical  
292 features within selected lipids. We applied the gradient-based interpretation method to the 813  
293 chosen molecular descriptors, assessing their contribution to the model's prediction. As illustrated  
294 in Figures 5a and 5b, we have visualized the top 20 salient descriptors for both the Hela cell line  
295 and RAW 264.7. For the Hela cell line, VSA\_EState3 and SssNH emerged as the most influential  
296 molecular characteristics for potency prediction. VSA\_EState3, a descriptor quantifying the  
297 electronic and steric properties of a molecule's surface area within a specific range<sup>50</sup>, along with  
298 SssNH, representing a tertiary amine, aligned with the expert understanding that head groups with  
299 tertiary amines are vital for lipid design. Subsequent analysis of essential features classified by  
300 head groups (Fig. 5k) pinpointed PEOE and Estate as the most critical descriptors for top-  
301 performing head groups (A13, A21), while SsNH2 (Sum of sNH2 E-states) and NsNH2 (Number  
302 of atoms of type sNH2) dominated in the least-performing groups (A5, A17) (Supplementary Fig.  
303 S4). Notably, these descriptors have strong associations with the amide bond in the structure, a  
304 critical connection within the 3CR Ugi Markush structure. This connection allows for various  
305 functional group attachments, influencing the lipid-like substances' overall charge and their  
306 physicochemical properties within biological systems. Intriguingly, the model does not favor  
307 amide bond generation, potentially due to its impact on the overall physicochemical properties of  
308 lipids, such as pKa. In the context of RAW 264.7, SpDiam\_Dzi and VR3\_D are identified as the  
309 most influential descriptors (Fig. 5b). VSA\_EState appears as the third most influential, implying  
310 its pivotal role in determining delivery potency to RAW 264.7, akin to Hela cells. Interestingly,  
311 head groups that underperformed in Hela (A5, A17) emerged as top performers in RAW 264.7,  
312 with SsNH2 and NsNH2 remaining the most critical features. In Hela cells, the cyclized head  
313 group outperformed the linear head group in transfection efficacy. However, the opposite trend  
314 was observed in RAW 264.7 cells. These observations underscore the necessity of designing LNPs  
315 with specific lipids tailored for distinct cellular targets.



317 Our subsequent analysis, as illustrated in Fig. 5e, explicates the relationships among lipid  
318 candidates targeting Hela cells, as identified by similarities in the AGILE model's lipid  
319 representations. We constructed a similarity network for the chosen 15 lipids, linking each lipid to  
320 its nearest equivalents. H9, the most potent LNP, demonstrated connections not only to LNPs with  
321 an identical head group (H7, H8) but also to other high-performing candidates, as identified by  
322 relative luciferase units (H12, H13). To gain further insights, we carried out molecular  
323 explanations on H9, illuminating the most salient regions in the molecule structure that heavily  
324 influenced the graph encoder's prediction within the AGILE model (Fig. 5c, Methods 1.8).  
325 Interestingly, head group structures emerged as the most salient for H9, which aligns with our  
326 previous findings emphasizing the importance of head groups. Similarly, we developed a similarity  
327 network for the 15 candidates selected for RAW 264.7 (Fig. 5f). R6 exhibited connections with  
328 other high-performing candidates, including R3, R8, and R11. These four lipids share identical tail  
329 structures: one being a C-12 alkyl chain and the other a C-18 alkyl chain. This shared characteristic  
330 suggests a strong correlation between these tail structures and the high transfection potency of R6,  
331 R3, R8, and R11. Interestingly, both tails are non-biodegradable, which hints at the potential  
332 necessity of lipid stability for successful macrophage transfection. Furthermore, these high-  
333 performing lipids commonly feature asymmetrical alkyl chains, a trait shared with SM102, which  
334 facilitates the formation of an inverted cone geometry more readily<sup>51</sup>. Similar to the findings for  
335 H9, head group structures were identified as an influential factor on the saliency map for R6.  
336 Additionally, the tail end was also highlighted as a salient region (Fig. 5d).

337 Moreover, our results highlight the importance of the carbon chain length of R2 as a critical factor  
338 in predicting transfection potency, particularly concerning RAW 264.7. It presents the distribution  
339 of predicted potencies relative to varying carbon chain lengths of R2 for lipids hailing from the  
340 top-performing head group A5 (Fig. 5g). Two distinct aspects emerge from this distribution: (1)  
341 As the R2 carbon chain length increases from 10 to 12, a corresponding uptick in predicted potency  
342 becomes apparent. Interestingly, any further extension in the R2 length inversely impacts the  
343 predicted potency. (2) In addition, R2's shorter carbon chain lengths ( $C \leq 12$ ) correlate with less  
344 variance in potency predictions compared to their longer counterparts ( $C > 12$ ). This trend is not  
345 restricted to the top-performing head groups but resonates across others, as well (Supplementary  
346 Fig. S14), a phenomenon further corroborated by a Pearson Correlation coefficient of -0.58.  
347 Examining the distribution plot for all lipids in the candidate set (Fig. 5h) reveals a similar pattern  
348 concerning R2 carbon chain length and predicted potency, albeit with a slightly attenuated Pearson  
349 Correlation of -0.39. Notably, we observe less variability amongst the shorter R2 chains ( $C \leq 12$ ).  
350 Interestingly, the importance of carbon chain lengths varies asymmetrically between the two  
351 respective tails. As shown in Fig. 5i, the correlation between the predicted potencies and R3 carbon  
352 chain lengths is noticeably lower than that of the R2 carbon chain lengths (-0.15 vs. -0.39). These  
353 tail-length findings pertain specifically to transfection in RAW 264.7. As displayed in Fig. 5j, the  
354 pattern within the Hela cell line is less defined, resulting in a Pearson correlation of -0.22.  
355 Collectively, these insights hold significant implications for guiding the design of LNPs  
356 specifically tailored for RAW 264.7.

357

358

## 359 Discussion

360 In this work, we presented the AGILE platform trained on comprehensive virtual and wet-lab  
361 libraries to enable predictions of LNP potency across different cell lines even in data-limited  
362 settings. Through exposure to an extensive array of molecular descriptors during the training  
363 process, the deep learning component in AGILE gained fundamental insights into the complex  
364 dynamics of LNP design, incorporating features like electronic and steric properties, and carbon  
365 chain lengths in a completely self-supervised manner.

366 One of the important findings is the influence of the molecular descriptor VSA\_EState3 and  
367 SssNH on the potency prediction for the HeLa cell line. These descriptors, which quantify the  
368 electronic and steric properties of a molecule's surface area and represent a tertiary amine,  
369 respectively, align with current expert understanding in lipid design. The connection between these  
370 molecular characteristics and their influence on delivery potency exemplifies the power of deep  
371 learning in elucidating nuanced molecular features. This correlation between expert knowledge  
372 and model interpretation endorses the validity of AGILE's predictive capabilities and lays a  
373 groundwork for future studies on other cell lines. Contrastingly, for the RAW 264.7 cell line,  
374 SpDiam\_Dzi and VR3\_D were identified as the most influential descriptors, highlighting the  
375 different physicochemical properties favored by different cell types. This variance in influential  
376 descriptors underscores the need for cell-specific LNP design, emphasizing the limitations in  
377 applying a one-size-fits-all approach to LNP design across diverse cell lines.

378 The molecular explanation applied to H9, the most potent LNP for the HeLa cell line, further  
379 corroborated the importance of head groups, a knowledge already prevalent in LNP design. On the  
380 other hand, for RAW 264.7, the high-performing LNPs shared identical tail structures, hinting at  
381 the potential role of tail structures in macrophage transfection. The fact that these tail structures  
382 are non-biodegradable also implies the significance of lipid stability in LNP potency. Such  
383 findings, which would be otherwise elusive without AGILE, elucidate the inherent complexities  
384 involved in tailoring LNPs for individual cell types.

385 Moreover, we found that the carbon chain length of R2 was a critical determinant of transfection  
386 potency, particularly in RAW 264.7. This result brings attention to the need for a delicate balance  
387 in the chain lengths to achieve optimal transfection, further complicating the LNP design process.  
388 The variance in the correlation between predicted potencies and carbon chain lengths for different  
389 tails - R2 and R3, as well as the asymmetric importance between the two respective tails, reinforces  
390 the idea that LNP design is a delicate process involving numerous factors and dependencies.

391 Furthermore, we found that AGILE's predictive power consistently improved with training on  
392 larger and more diverse datasets, mirroring observations in fields like natural language  
393 understanding, computer vision, and mathematical problem-solving. The exposure to extensive  
394 datasets during training also seemed to enhance AGILE's robustness to various factors and  
395 dependencies involved in LNP design. These findings suggest that as we continue to expand our  
396 dataset, future models pretrained on even larger scales may yield more precise predictions in  
397 elusive tasks with increasingly limited task-specific data. For example, beyond using AGILE to  
398 discover LNPs for mRNA delivery to previously unexplored tissues and cell types, there's an

399 opportunity to expand the wet-lab mRNA transfection data from cell cultures to in vivo data from  
400 animal studies and ex vivo data in human tissues. This could potentially boost the efficiency and  
401 reliability of LNPs discovered by AGILE for in vivo mRNA delivery in human patients, thereby  
402 supporting the clinical development of mRNA LNP products. Additionally, incorporating more  
403 diverse combinatorial chemistry methods, along with comprehensive wet-lab data, could further  
404 enhance the chemical diversity for AGILE model training<sup>52</sup>. This could allow AGILE to identify  
405 ionizable lipids with specific functionalities, such as immunostimulatory properties, essential for  
406 mRNA vaccine delivery and cancer immunotherapy. Furthermore, AGILE could adopt recent  
407 generative models, like diffusion networks<sup>53, 54</sup>, to generate *de novo* lipid molecules for specific  
408 applications.

409 Overall, AGILE synergizes the strength of combinatorial chemistry and deep learning, elucidating  
410 the intricate dynamics of LNP design and making this insight accessible for a multitude of  
411 downstream applications. AGILE's ability to identify and interpret influential molecular  
412 descriptors represents a significant leap forward in the field of nanomedicine, particularly in lipid  
413 design. Its capacity in predicting the transfection efficacy of LNPs in diverse cell lines, including  
414 challenging ones like macrophages, holds promise for not only improving mRNA delivery but also  
415 for guiding CAR cell therapy and other immunotherapeutic strategies. It can potentially accelerate  
416 the discovery of potent LNPs and facilitate the design of tailored ionizable lipids for mRNA  
417 delivery, thereby contributing significantly to the continuous development of mRNA-based  
418 therapeutics and their deployment in clinical settings.

## 419 **Materials and methods**

420 Extended materials and methods are available in the supplementary information (SI).

### 421 **1.1 Data Preparation**

#### 422 **Virtual Library**

423 We utilized Ugi combinatorial chemistry method to design diverse head groups, connecting groups,  
424 and two distinct alkyl chains. To be specific, we used the Markush Editor in the ChemAxon Marvin  
425 Suite (Marvin 23.4.0, ChemAxon, <https://www.chemaxon.com>). The resulting virtual library  
426 contained approximately 60,000 lipid structures which were then exported into SMILES strings.  
427 This virtual library comprises multiple carbon chains, from C6 to C26. In addition, the presence  
428 or absence of ester bonds and their position in the carbon chain are used to improve the chemical  
429 diversity of the virtual library. The surface charge of LNP is usually determined by the lipids' head  
430 groups. In addition, the head group is critical for mRNA binding. Amine groups are commonly  
431 used as lipids' head groups to form hydrogen bonds with mRNA, especially those containing  
432 tertiary amine.

#### 433 **Experimental library**

434 Our experimental library contains 20 head groups, 12 carbon chains with ester bonds, and 5 carbon  
435 chains with isocyanide head groups. We selected 1200 lipids for chemical synthesis and *in vitro*  
436 transfection potency experiments in Hela and RAW 264.7 cell lines. We label the corresponding

437 mRNA transfection potency in cells to each compound for the 1,200 lipids library. And these data  
438 are generated by ChemAxon Marvin Suite into SMILE files (SMILE files in SI).

### 439 **Candidate library**

440 The final library used for model prediction is a filtered subset of the virtual library. The filtering  
441 contains three steps based on availability and rationality. First, we retained the lipids containing  
442 tertiary amine structures. Second, we removed tail chains that were too long (>C18) or too short  
443 (<C10) based on expert knowledge of plausible ionizable lipid design<sup>36</sup>. Last, we select only those  
444 reagents commercially available for further validation of the model. Upon completion of the  
445 filtering process, the final candidate library comprises approximately 12,000 lipids (SMILE files  
446 in SI), with 22 unique head groups (Supplementary Fig. S17), and a distinct arrangement of 9 R2  
447 tail types alongside 2 R3 tail types (Supplementary Fig. S18). In the prediction step of the platform,  
448 the model proposed the most promising lipids by predicting and ranking on the candidate library.

### 449 **1.2 Molecular graph construction**

450 Molecular structures can be naturally represented as graphs where atoms are nodes and bonds are  
451 edges. For each molecule, the SMILES representation is converted into a molecular graph using  
452 RDKit<sup>55</sup>, and later input to the neural network model in the platform. This representation captures  
453 the topological structure and properties of a molecule effectively. An LNP molecule graph  $G$  is  
454 defined as  $G = (V, E)$ , where nodes  $V$  represent the atoms and edges  $E$  represent chemical bonds.  
455 The atom node features include the atom type (as on the periodic table) and a flag indicating  
456 whether the whole molecule it belongs to is chiral. For a node  $v$ , the features are constructed in a  
457 two-dimensional vector,  $h_v \in N^2$ . Edge features are constructed based on respective chemical  
458 bond types (i.e., single, double, triple, or aromatic bonds) and the stereochemical directionality  
459 (i.e., the `rdchem.BondDir` in RDKit). Similarly, the edge features form another two-dimensional  
460 vector for each bond between atom  $v$  and  $u$ ,  $\epsilon_{v,u} \in N^2$ .

### 461 **1.3 The Model Architecture**

462 The deep learning model in AGILE comprises three major components: (1) The embedding layers  
463 to project node and edge features into learnable vectors, (2) the graph encoder for modeling  
464 molecular structures, and (3) the descriptor encoder for modeling molecular properties.

#### 465 **Embedding Layers**

466 The embeddings layers project the integer features in  $h_v$  and  $\epsilon_{v,u}$  to learnable feature vectors  
467  $h_v^{(0)}$  and  $\epsilon_{v,u}^{(0)}$ , which can be optimized later during the training of the whole neural network. Here,  
468 both  $h_v^{(0)}$  and  $\epsilon_{v,u}^{(0)}$  are  $R^d$  vectors, and  $d$  is a predefined size of embedding dimensions. To be  
469 specific, we first obtained the embedding vectors for both atom type and charity features in  $h_v$ ,  
470 and added the two vectors elementwise to output the  $h_v^{(0)}$ :

$$h_v^{(0)} = Emb_{h,0}^{(0)}(h_v[0]) + Emb_{h,1}^{(0)}(h_v[1]), \quad Eq. 1$$

471 here [i] denotes the i-th element in the vector. *Emb* is the embedding layer projection. In this work,  
472 we use the PyTorch Embedding layers ([https://pytorch.org/docs/stable/generated/](https://pytorch.org/docs/stable/generated/torch.nn.Embedding.html)  
473 [torch.nn.Embedding.html](https://pytorch.org/docs/stable/generated/torch.nn.Embedding.html)). Similarly, the  $\epsilon_{v,u}^{(0)}$  is computed as:

$$\epsilon_{v,u}^{(0)} = Emb_{\epsilon,0}^{(0)}(\epsilon_{v,u}[0]) + Emb_{\epsilon,1}^{(0)}(\epsilon_{v,u}[1]). \quad Eq. 2$$

474

## 475 Graph Encoder

476 We used Graph Isomorphism Network (GIN)<sup>56</sup>, a type of graph neural network (GNN), to operate  
477 on the input molecule graphs and to learn a representation vector for each LNP molecule. GIN can  
478 directly propagate messages among nodes and edges on a graph structure and thus is suitable for  
479 processing molecular graphs. Additionally, the advantage of GIN over other GNNs is its ability to  
480 distinguish between different graph structures, including isomorphic graphs. This makes GIN more  
481 expressive than many other GNNs and a suitable tool for tasks involving molecular graph data. It  
482 is worth noting that the implemented GIN model follows the similar structures used in MolCLR<sup>57</sup>,  
483 so that we can benefit from the general pretrained molecular model of MolCLR as a warm start  
484 for the platform (Methods 1.4). The update rule of GIN for a node representation on the  $k^{th}$  layer  
485 is given as:

$$h_v^{(k)} = MLP^{(k)} \left( (1 + \epsilon^{(k)}) \cdot h_v^{(k-1)} + \sum_{u \in N(v)} m_u^{(k-1)} \right), \quad Eq. 3$$

486 where  $h_v^{(k)}$  is the representation of node  $v$  at the  $k^{th}$  layer and  $N(v)$  denotes the set of neighbors  
487 of node  $v$ , and  $\epsilon$  is a learnable parameter. MLP denotes the stacked fully connected neural network  
488 layers. The  $m_u^{(k-1)}$  is the message propagated between a neighbor  $u$  to the current node. It is  
489 computed as the sum of node and edge contributions:

$$m_u^{(k-1)} = h_u^{(k-1)} + \epsilon_{v,u}^{(k-1)}, \quad Eq. 4$$
$$\epsilon_{v,u}^{(k-1)} = Emb_{\epsilon,0}^{(k-1)}(\epsilon_{v,u}[0]) + Emb_{\epsilon,1}^{(k-1)}(\epsilon_{v,u}[1]).$$

490 Notably, we use  $h_v^{(0)}$  and  $\epsilon_{v,u}^{(0)}$  from Eq. 1 and Eq. 2 for the first GIN layer.

491 We stack a total of  $K$  GIN layers for the entire Graph Encoder. To extract the feature of the whole  
492 molecular graph  $h_G$ , we implemented the mean pooling operation on the final layer to integrate all  
493 the node features:

$$h_G = Mean(\{h_v^{(K)} : v \in G\}). \quad Eq. 5$$

494 Another fully connected layer is used to transform  $h_G$  to the final lipid representation  $z_G$ :

$$z_G = MLP(h_G). \quad Eq. 6$$

495

## 496 Molecular Descriptor Encoder

497 In addition to the structure features encoded by the GIN, the platform utilizes another descriptor  
498 encoder to explicitly model molecular properties. In our experiment, we found this contributes a  
499 more stabilized training optimization. We hypothesize that this benefit come from the straight-  
500 forward utilization of computed properties during the optimization, which relieves the model from

501 learning all information from the structure alone. In the implementation of the platform, the  
502 molecular descriptors derived from Mordred<sup>30</sup> calculations were used, which contain over 1,000  
503 common descriptors for each molecule, including the num of atoms, num of bonds, et. al. These  
504 features are encoded by gully connected layers into a representation for these properties,  $z_p \in R^{d_p}$ :

$$z_p = MLP(descriptors). \quad Eq. 7$$

505 The final representation of the molecule is the concatenation of the structure and property  
506 representations:

$$z = [z_G, z_p], \quad Eq. 8$$

507 where  $[ , ]$  denotes the concatenation of two vectors.

## 508 1.4 Model Pre-training

509 The model pre-training aims to learn generalizable lipid representation that can benefit the  
510 downstream transfection potency prediction task. Before our lipid-oriented pre-training, we first  
511 initialized the model parameters by the general pre-trained model from MolCLR<sup>57</sup>, which has been  
512 trained on over 10 million distinct small molecules. The rationale for this initialization is to provide  
513 a warm start of a model that already has been trained to capture molecular structures. Next, we  
514 perform continuous pre-training on the 60,000 lipids in the virtual library (Methods 1.1) using  
515 contrastive learning to optimize the model's performance within the LNP domain.

### 516 Contrastive learning objective

517 Our pre-training objective is to learn LNP representation through contrasting positive data pairs  
518 against negative pairs. The model is trained to minimize the following loss:

$$L_{i,j} = -\log \frac{\exp\left(\frac{\text{sim}(z_i, z_j)}{\tau}\right)}{\sum_{k=1}^{2N} \mathbb{I}\{k \neq i\} \exp\left(\frac{\text{sim}(z_i, z_k)}{\tau}\right)}, \quad Eq. 9$$
$$\text{sim}(z_i, z_j) = \frac{z_i z_j}{\|z_i\|_2 \|z_j\|_2},$$

519 where  $z_i$  and  $z_j$  are the learned lipid representation vectors extracted from a positive data pair,  $N$   
520 is the batch size, and  $\tau$  is the temperature parameter set manually. In this pre-training step, we  
521 omitted the descriptor encoder, so the lipid representation only contains the graph structure  
522 representation  $z_G$  as in Eq. 6. To construct the positive data pair, each input lipid molecule graph  
523 is transformed into two different but correlated molecule graphs using graph augmentation. The  
524 molecule graphs augmented from the same molecule are denoted as a positive pair, and those from  
525 different molecules are denoted as negative pairs within each batch. During training, the model  
526 learns to maximize the agreement of positive pairs while minimizing the agreement of negative  
527 ones.

### 528 Data Augmentation

529 We used two augmentation strategies inherited from the MolCLR<sup>57</sup> pre-training workflow at the  
530 atom and bond levels. In the continuous pre-training of LNP molecules, three molecular graph data  
531 augmentation strategies are consistently employed. 1) Atom masking: Within the lipid molecular  
532 graph, atoms are randomly masked according to a specified ratio. This process compels the model  
533 to assimilate chemical information, such as atom types and corresponding chemical bond varieties  
534 within lipid molecules. 2) Bond deletion: Chemical bonds interconnecting atoms are randomly  
535 removed in accordance with a designated ratio. As the formation and dissociation of chemical  
536 bonds dictate the properties of LNP molecules during chemical reactions, bond deletion facilitates  
537 the model's learning of correlations between LNP molecule involvement in various reactions.

## 538 **1.5 Model Fine-tuning**

539 The lipid-oriented pretrained model (Methods 1.4) serves as the starting point of the fine-tuning  
540 stage. During the fine-tuning, we included the Molecular Descriptor Encoder and used the  
541 combined output  $z$  in Eq. 8 as the molecule representation. For the property descriptor input, a  
542 series of preprocessing procedures are executed, aiming to isolate pertinent features. Initially,  
543 descriptors with a standard deviation of zero are eliminated, followed by the selection of  
544 descriptors exhibiting correlation with the experimentally determined transfection potency in both  
545 HeLa and Raw 264.7 cells (R2 score > 0.006), resulting in the identification of 813 salient  
546 descriptors (Supplementary Fig. S19). Subsequently, log transformation is applied to descriptors  
547 possessing extensive data ranges, with normalization conducted accordingly. The preprocessing  
548 steps enacted on the fine-tuning dataset are documented and replicated for the 12,000 lipids in the  
549 candidate library in anticipation of the model prediction phase (Methods 1.6).

550 The model is fine-tuned utilizing the 1,200 lipids of the experiment library to perform regression  
551 on LNP transfection potency. The mean squared loss between the predicted and ground-truth  
552 potency is used to optimize the model parameters:

$$L_{mse} = \frac{1}{n} \sum_{i=1}^n (Pred(z_i) - y_i)^2, \quad Eq. 10$$

553 where  $Pred(\cdot)$  denotes the fully connected layers that perform the potency prediction, and  $y_i$  is  
554 the actual transfection potency recorded *in vitro*.

555 A scaffold-based 80%-10%-10% train-valid-test split is performed on the experimental library.  
556 We fine-tune the model on the training set only and evaluate the performance on the validation set  
557 using root mean squared error (RMSE) and Pearson correlation with the ground truth transfection  
558 potency.

## 559 **1.6 Model ensemble prediction and candidate ranking**

560 To enhance the model's robustness and generalizability, the fine-tuning process is carried out ten  
561 times, from which the top five models are selected based on RMSE and Pearson correlation  
562 performance on the testing set. These five models are subsequently employed for ensemble  
563 prediction on the 12,000-member candidate set. We first get the potency predictions from each  
564 model and calculate the average and standard deviation of the five predicted values for each  
565 candidate molecule. The mean predicted values are then subtracted from the standard deviation,  
566 and the resulting predicted score is used to rank the candidates.  
567 We observed that the predicted potencies exhibit distinct stratification by head groups and tail  
568 combinations, and the structural differences among molecules with the same head groups and tail  
569 combinations are relatively minor (Supplementary Fig. S20). To increase the diversity of selected  
570 candidates, we implement a ranking scheme that sorts candidate LNPs by head groups and tail  
571 combinations (Supplementary Fig. S21). Given the predicted values, candidates are first organized  
572 by head groups and subsequently ranked in descending order. Candidates within each head group  
573 are then ranked by tail combinations following the same schema. Ultimately, we select the top five  
574 head groups and the top three tail combinations from each head group, resulting in a final candidate  
575 set of 15 LNPs.

## 576 **1.7 Implementation details**

577 The Graph Encoder in the model consists of a five-layer GIN with ReLU activation. To extract a  
578 512-dimensional LNP representation, an average pooling layer is applied to each lipid molecular  
579 graph. A single hidden layer MLP is then employed to map the representation into a 256-  
580 dimensional latent space. During model pre-training, the contrastive loss is optimized using the  
581 Adam optimizer<sup>58</sup>, with a weight decay of  $10^{-5}$ , and the temperature is set to 0.1. The pre-training  
582 process involves a batch size of 512 for 100 epochs.

583 For model fine-tuning, an additional MLP with one hidden layer is introduced to map the molecular  
584 descriptors into 100-dimensional latent vectors. These vectors are concatenated with the 256-  
585 dimensional LNP representation obtained from the GNN encoder. Subsequently, a two-layer MLP  
586 is utilized to derive the final prediction value from the concatenated vector. The fine-tuning process  
587 employs the Adam optimizer with a weight decay of  $10^{-6}$  to optimize the loss (Eq. 10). Each fine-  
588 tuned model is trained using a batch size of 128 for 30 epochs.

## 589 **1.8 Model interpretation**

### 590 **Salient molecular descriptors calculation**

591 In our study, we employed the Integrated Gradients<sup>59</sup> methodology featured in the Captum<sup>60</sup>  
592 Python package to interpret the significance of molecular descriptors. The process involves  
593 approximating the integral of molecular descriptor gradients in relation to their respective  
594 predicted potencies for each LNP within the candidate library. A molecular descriptor's  
595 prominence is proportionate to the absolute value of its integrated gradient. We implemented  
596 computations across all five ensemble models for each target cell line. To calculate an overall  
597 significance for each feature, we initially averaged the computed gradients across all input samples



598 on each model, subsequently normalizing these importance scores. The final step involved  
599 computing the mean of these importance scores across all five models. The top 20 critical features  
600 were selected and visualized based on the calculated importance scores. When assessing feature  
601 significance in the context of head groups, we averaged the integrated gradients for each head  
602 group and then proceeded to normalization. Following this, we averaged the results across the five  
603 models for each respective head group. The top two significant features for each head group were  
604 then selected, and their scores were visualized across all head groups.

### 605 **Construction of the similarity network on the selected candidates**

606 We constructed a similarity network for the 15 selected candidates respective to each target cell  
607 line, with the aim of elucidating the similarities among the candidates. Utilizing the LNP vector  
608 representations provided by the corresponding fine-tuned model, we computed the cosine  
609 similarities for each candidate pair and chose the four most similar neighbors for each. This  
610 generated similarity network was then visualized, with the node sizes representing the relative  
611 luciferase units.

### 612 **Molecular structure interpretation**

613 To ascertain the critical areas within the LNP structure that contribute significantly to the model's  
614 predictions, we engaged the Model Agnostic Counterfactual Compounds Generation feature  
615 present in the ExMol Python package<sup>61</sup>. This is accomplished by generating molecular  
616 counterfactuals and investigating the alterations required in the LNP molecule to modify its  
617 predicted transfection potency (Supplementary Fig. S22). The molecular counterfactuals produced  
618 are designed to retain as much similarity to the input LNP molecule as feasible. If modifications  
619 in particular regions result in either an increase or decrease in the predicted potency, such areas  
620 are deemed as essential regions. The critical areas identified through this process were visualized  
621 for both H9 and R6.

## 622 **1.9 Materials and lipid library synthesis**

623 To prepare our materials, we got amines and starting compounds from Sigma-Aldrich and TCI  
624 America. We then put 10  $\mu$ L of a 350  $\mu$ M stock solution containing amines and tails into each well  
625 of a 96-well plate with glass inserts. This stock solution was made by mixing the compounds in a  
626 2:1 ratio of methanol with 0.2 eqv. catalyst phenyl hypophosphoric acid ( $H_3PO_4$ ). The plates were  
627 covered and placed on a shaker to stir overnight, with conversions yield typically over 70%. We  
628 also formulated lipids into LNP in the same reaction plates. These lipids were purified through  
629 flash column chromatography, and their final structures were confirmed using  $^1H$  400 MHz NMR  
630 spectrometry with  $CDCl_3$  and tetramethylsilane (TMS) as a standard at UHN Nuclear Magnetic  
631 Resonance Core Facility. To further analyze our materials, we obtained high-resolution mass  
632 spectra using an LC-Mass spectrophotometer at the Centre for Pharmaceutical Oncology of the  
633 University of Toronto.

## 634 **1.10 LNP synthesis and formulation for high throughput screening**

635 To conduct high-throughput screening, we created an organic phase by dissolving a mixture of  
636 cationic lipid, DOPE (Avanti), cholesterol (Chol, Sigma-Aldrich), and C14-PEG 2000 (Avanti) in  
637 ethanol at a predetermined molar ratio. We prepared the aqueous phase using firefly luciferase  
638 mRNA (mLuc, Translate), Cre recombinase mRNA (TriLink BioTechnologies) or EGFP-mRNA  
639 (TriLink BioTechnologies) in 10 mM sodium citrate buffer (pH 4.0, Fisher). All mRNAs were  
640 stored at -80 °C and were allowed to thaw on ice before use. During the high-throughput screening  
641 phase, LNPs were synthesized by mixing an aqueous phase containing the mRNA with an ethanol  
642 phase containing the lipids by the OT-2 pipetting robot. The aqueous phase was prepared in a 10  
643 mM citrate buffer with the corresponding mRNA. The ethanol phase was prepared by solubilizing  
644 a mixture of ionizable lipid, helper phospholipid (DOTAP, DOPE, cholesterol, and C14-PEG 2000  
645 at pre-determined molar ratios with an ionizable lipid/mRNA weight ratio of 10 to 1.

### 646 **1.11 LNP synthesis and formulation for in vitro and vivo tests**

647 For other *in vitro* and *in vivo* tests, all materials were prepared and processed without nucleases  
648 throughout the synthesis and formulation steps. DLin-MC3-DMA and ALC0315 were purchased  
649 from Echelon Biosciences. MC3-LNP was prepared at the molar ratio of 50:10:38.5:1.5 (DLin-  
650 MC3-DMA:DSPC: cholesterol: DMG-PEG2000) and ALC0315-LNP was prepared at the molar  
651 ratio of 46.3:9.4:42.7:1.6 (ALC0315:DSPC: cholesterol: ALC0159 [Echelon Biosciences]). The  
652 optimal formulations of H278 and R080 LNPs for the subsequent experiments were determined  
653 by the LNP formulation optimization method. Except for the high-throughput screening, the  
654 aqueous and ethanol phases were rapidly mixed by pipette at a 3:1 volumetric ratio. Post incubation  
655 for 15 min in a 4 °C fridge.

### 656 **1.12 LNP formulation optimization.**

657 The statistical software JMP 16 (SAS Institute) analyzed the experimental data. In this Design of  
658 experiments (DoE) approach, the four-factor Box-Behnken design was suitable for second-order  
659 models comprising 17 preparation runs. The design was cited as a common experimental design  
660 for screening crucial factors. In this design, all factors (lipid/mRNA weight ratio, ionizable lipid  
661 molar ratio, helper lipid molar ratio, and PEG molar ratio) have low, center, and high levels.

### 662 **1.13 In vitro high throughput screening.**

663 The lipid library, which was not purified, was directly combined with ethanol and the aqueous  
664 solution of mLuc. For *in vitro* transfection, the lipid-mRNA mixture, containing 0.1 µg of mRNA,  
665 was added to pre-seeded HeLa and Raw264.7 cells in 96-well plates. Following overnight  
666 incubation, the transfection potency of mLuc was measured using the One-Glo Luciferase Assay  
667 System (Promega), following the manufacturer's instructions. The luminescence was quantified  
668 using the Cytation imaging reader (BioTek). Finally, the resulting bioluminescence values are  
669 assigned to each SMILE string.

### 670 **1.14 In vivo luciferase mRNA for bioluminescence.**

671 At 6 h after the intramuscular administration of the mRNA LNPs, mice were injected  
672 intraperitoneally with 0.2 ml d-luciferin (10 mg/ml in PBS). The mice were anesthetized in a  
673 ventilated anesthesia chamber with 1.5% isoflurane in oxygen and imaged 10 min after the  
674 injection with an *in vivo* imaging system (IVIS, PerkinElmer). Luminescence was quantified using  
675 the Living Image software (PerkinElmer). C57BL/6 mice (4-8 weeks) were purchased from the  
676 Jackson Laboratories.

### 677 **1.15 ROSA<sup>mT/mG</sup> Cre reporter mice transfection analysis.**

678 All animal studies were approved and conducted in compliance with the University Health  
679 Network Animal Resources Centre guidelines. For gene recombinant Cre mRNA delivery, LNPs  
680 co-formulated with Cre mRNA (0.5 mg kg<sup>-1</sup>) were i.m. injected into ROSA<sup>mT/mG</sup> Cre reporter mice  
681 (The Jackson Laboratory). After 7 d, mice were killed, and major organs were collected and  
682 imaged using an IVIS imaging system (PerkinElmer). For direct fluorescence imaging, organs and  
683 muscle tissues were fixed in 4% buffered paraformaldehyde overnight at 4°C, then equilibrated in  
684 30% sucrose overnight at 4°C before freezing in OCT. Three nonconsecutive sections from each  
685 organ sample were mounted with DAPI to visualize nuclei and imaged for DAPI, tdTomato, and  
686 GFP. Sectioned into 10 µm depth, and further imaged using a Fluorescence microscope (Zeiss  
687 AXIO Observer 7 Inverted LED Fluorescence Motorized Microscope).

### 688 **1.16 Intracellular delivery of GFP mRNA to RAW 264.7**

689 For GFP mRNA delivery, GFP mRNA LNPs containing 500 ng GFP-mRNA were added to 24-  
690 well plates for 48 h incubation at 37 °C. Finally, a fluorescence microscope (Zeiss AXIO Observer  
691 7 Inverted LED Fluorescence Motorized Microscope) was used to evaluate the transfection effect.

### 692 **1.17 Statistical analysis**

693 The data were subjected to statistical analyses using GraphPad Prism 9 (GraphPad Software). A  
694 two-tailed unpaired Student's t-test was conducted to assess the significance of the comparisons as  
695 indicated. Data are expressed as mean ± s.d. P values <0.05 (\*), P < 0.01 (\*\*), P < 0.001 (\*\*\*) and  
696 P < 0.0001 (\*\*\*\*) were statistically significant.

### 697 **Acknowledgments**

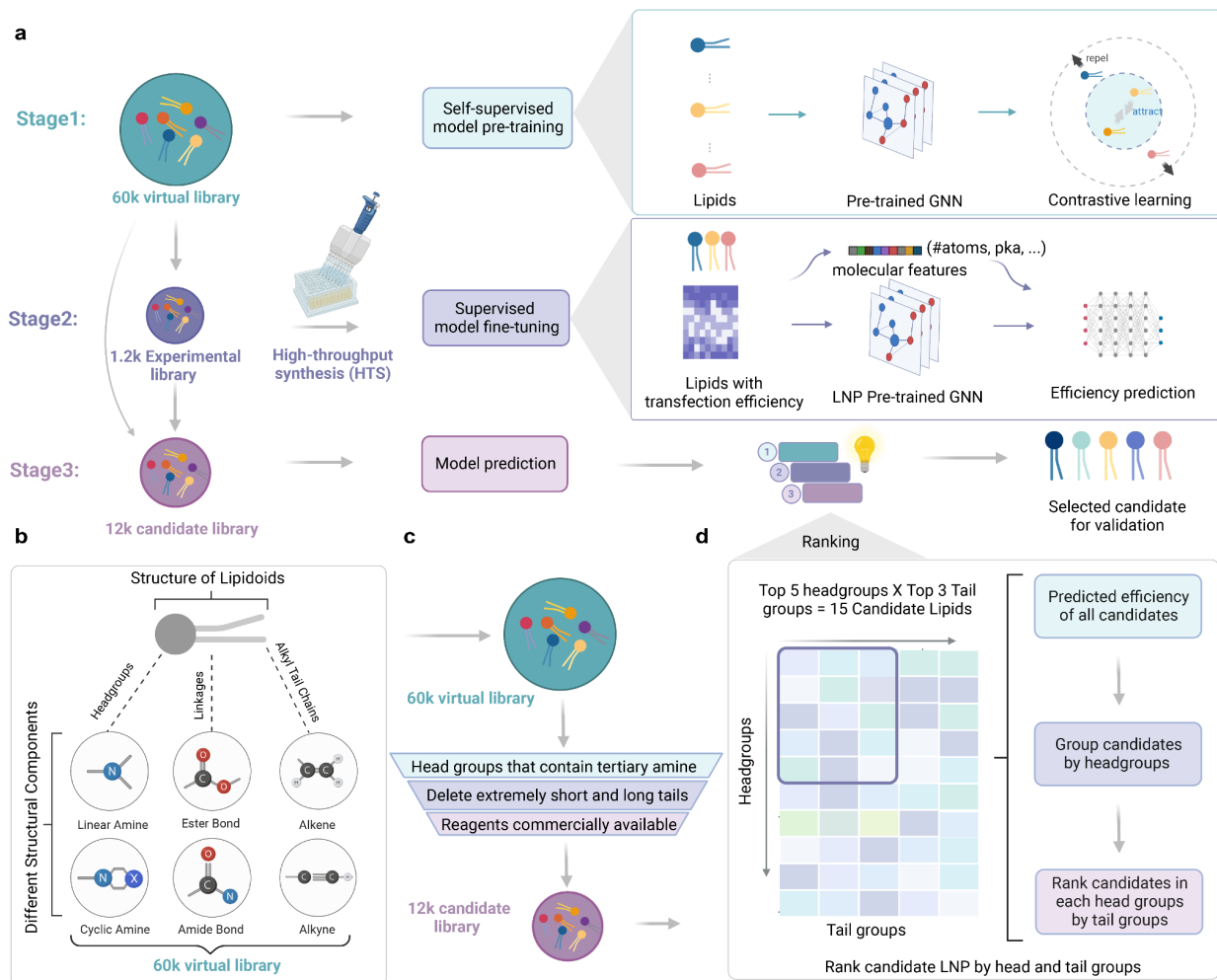
698 The authors are grateful to R.S. Langer for project discussions and constructive input. This work  
699 was supported by the Leslie Dan Faculty of Pharmacy startup fund, the Princess Margaret Cancer  
700 Center operating fund, the Connaught Fund (no. 514681), the J. P. Bickell Foundation (no.  
701 515159), the Canada Research Chairs Program (no. CRC-2022-00575), Canadian Institutes of  
702 Health Research (no. PJH-185722), Natural Sciences and Engineering Research Council of  
703 Canada (no. RGPIN-2023-05124) and the Canada Foundation for Innovation - John R. Evans  
704 Leaders Fund (no. 43711); Y.X. acknowledge the Postdoctoral Fellowship from PRiME-UHN  
705 Clinical Catalyst Program (no. PRMUHN2022-005); A.V. acknowledges the Postdoctoral  
706 Fellowship from the PRiME - Precision Medicine initiative at the University of Toronto; R.X.Z.L.  
707 acknowledges the Postdoctoral Fellowship from the Acceleration Consortium at the University of

708 Toronto. The authors acknowledge the technical support from the Centre for Pharmaceutical  
709 Oncology in Flow Cytometry, and Imaging Facilities, and acknowledge the Princess Margaret  
710 Cancer Centre for the use of NMR and Animal facilities. Balloon plots created with  
711 bioinformatics.com.cn. Figures 1-4 were created with Biorender.com.

712 **Competing Interests**

713 Y.X., J.C., and B.L. have filed a provisional patent for the development of the described lipids.

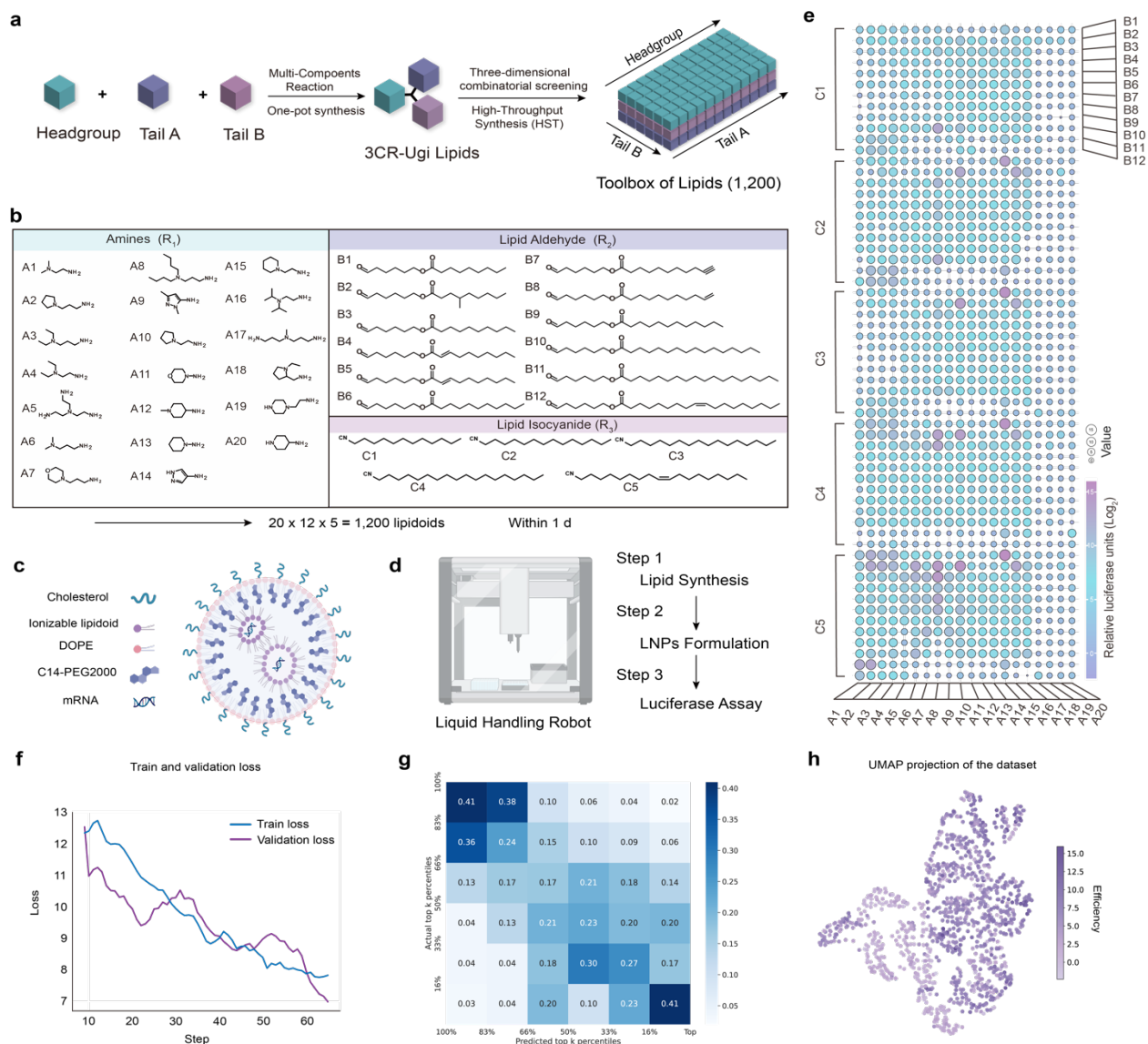
714



715

716 **Figure 1. Overview of the platform design pipeline.**

717 (a) Illustration of the 3-stage workflow of the platform. Stage 1: Construction of a virtual library  
718 and self-supervised pre-training of the model. Stage 2: Synthesis of an experimental library for the  
719 fine-tuning of the model in a supervised manner. Stage 3: Deployment of the fine-tuned model for  
720 predictive analysis on a candidate library, followed by ranking for final candidate selection. (b)  
721 Depiction of virtual library design through the application of Ugi combinatorial chemistry. (c)  
722 Schematic representation of the rational selection process for lipid candidates, with 3 listed  
723 filtering criteria. (d) A comprehensive breakdown of the ranking procedure and the selection  
724 methodology for final candidates.

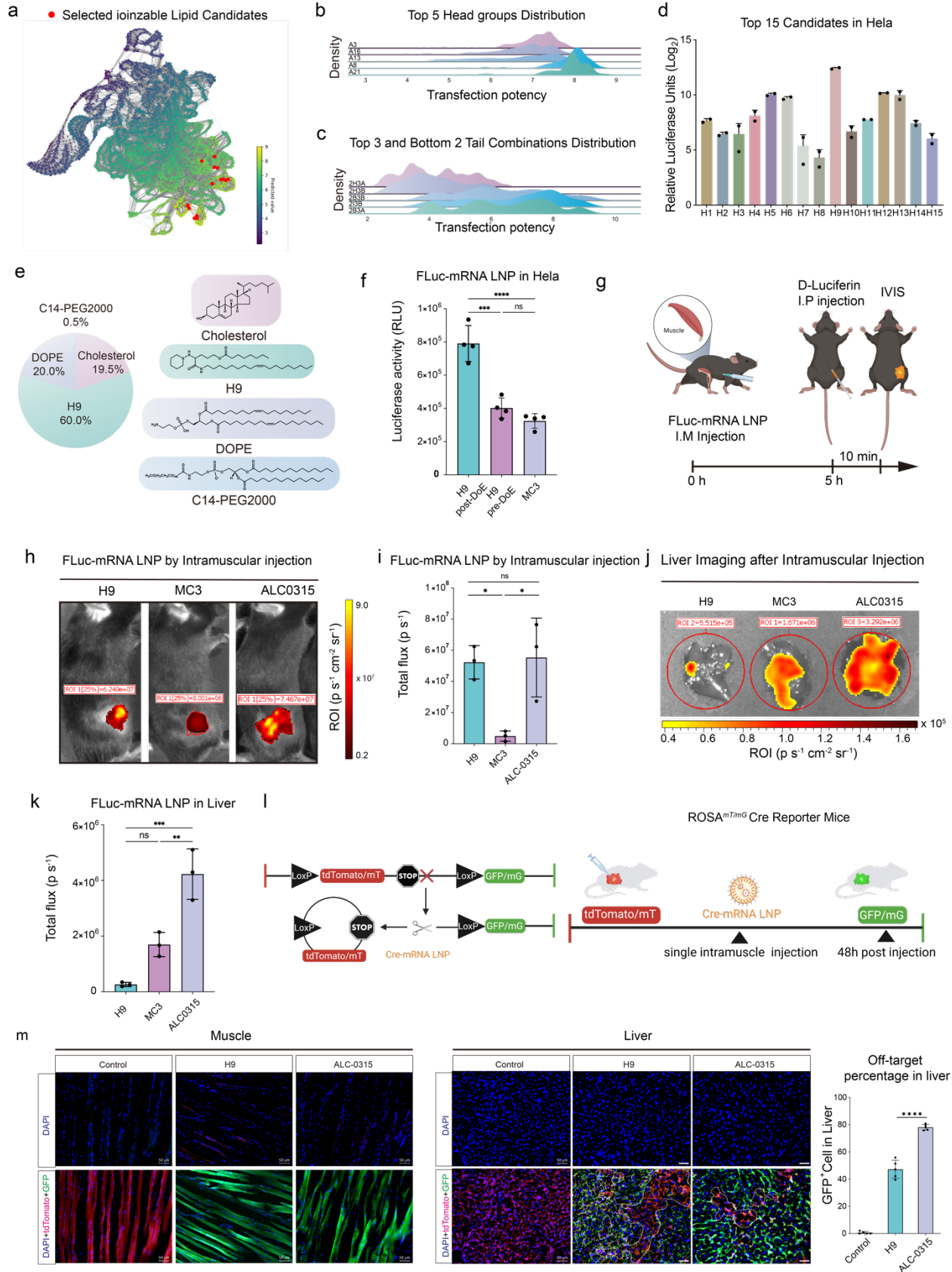


725

726 **Figure 2. High throughput lipids synthesis and screening platform.**

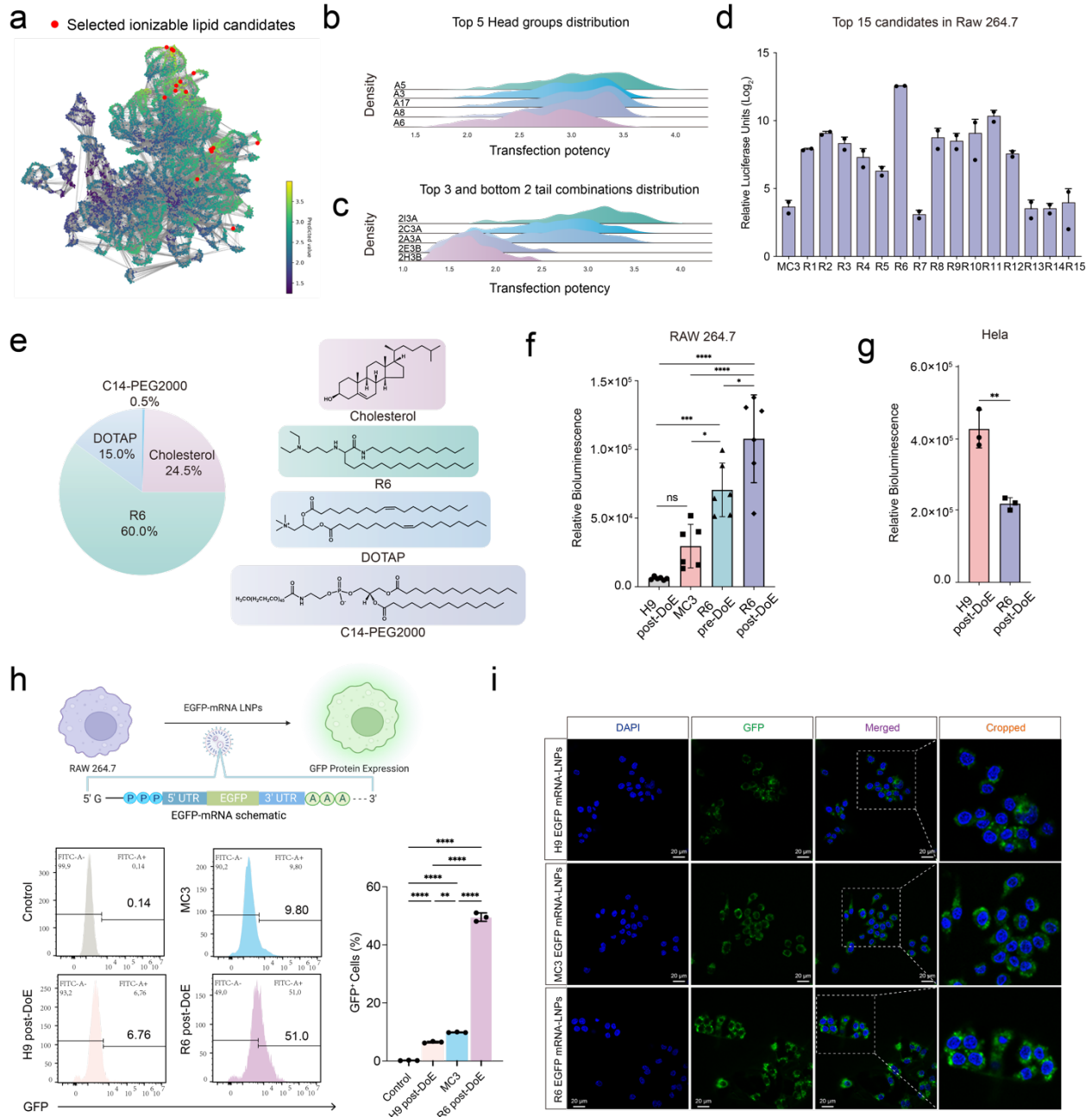
727 (a) A schematic to illustrate the high-throughput synthesis method for lipids. (b) The combinatorial  
 728 lipoids library consists of three components structure (amine head groups, aldehyde tails, and  
 729 isocyanide tails). (c) A schematic diagram shows the LNPs components for mRNA encapsulating.  
 730 (d) Lipid synthesis, LNPs formulation and luciferase assay based on liquid handling robot. (e) The  
 731 data used for the fine-tuning are depicted in a balloon plot, which involved 1,200 LNPs for Fluc  
 732 mRNA (mLuc) delivery and measuring the relative luciferase expression in Hela cells. (f) The loss  
 733 value on the training set and validation set against fine-tuning steps. (g) The precision matrix  
 734 computed on the experimental library of 1,200 lipids. The predicted and actual transfection  
 735 potencies are divided into six equal percentiles. (h) UMAP plot of the experimental library, colored  
 736 by the transfection potency.

737



739 **Figure 3. Model prediction and the validation of the gene editing potential with top-**  
740 **performing mRNA-LNPs.** (a) The UMAP plot of the predicted molecule trans potencies. (b)  
741 Head group distribution and (c) tail combination distribution in Hela. (d) Validate 15 lipid  
742 candidates for Hela cell. (e) The top-performing formulation parameters used in the optimization  
743 of H9 LNPs in Hela. (f) Transfection of mFFL LNPs in Hela cells (n = 4 biologically independent  
744 experiments per group). (g) A schematic to illustrate the Intramuscular (IM) injection of mFluc-  
745 loaded LNPs into the mice and IVIS imaging. (h) LNPs formulated with FFL encoding mRNA  
746 were injected intramuscularly into mice (0.25mg mRNA/kg mouse). The top-performing lipids H9  
747 with optimized formulation compare with the MC3 and ALC-0315 LNPs (n = 3 biologically  
748 independent mice per group, 0.5 mg kg<sup>-1</sup> mLuc per mouse). (i) Transfection of mFFL LNPs at the  
749 i.m. injection site in mice (n = 3 biologically independent mice per group). (j) IVIS imaging for  
750 liver after IM injection of mFluc-loaded LNPs. (k) Transfection of mFFL LNPs of liver in mice  
751 after IM injection (n = 3 biologically independent mice per group). (l) A schematic illustrating the  
752 Cre recombinase deletes STOP cassettes and activates the GFP mice reporter. (M) Representative  
753 confocal microscopy images and quantification of tdTomato and GFP expression in histological  
754 muscle and liver sections of mTmG mice post-injection of Cre-mRNA loaded LNPs by IM  
755 injection. Scale bar: 50µm. n = 5 sections from 3 mice. Error bars are S.D. Statistical significance  
756 was analyzed by the two-tailed Student's t-test. \* = p-value < 0.05, \*\* = p-value < 0.01, \*\*\* = p-  
757 value < 0.005. Data are presented as mean ± SD.  
758

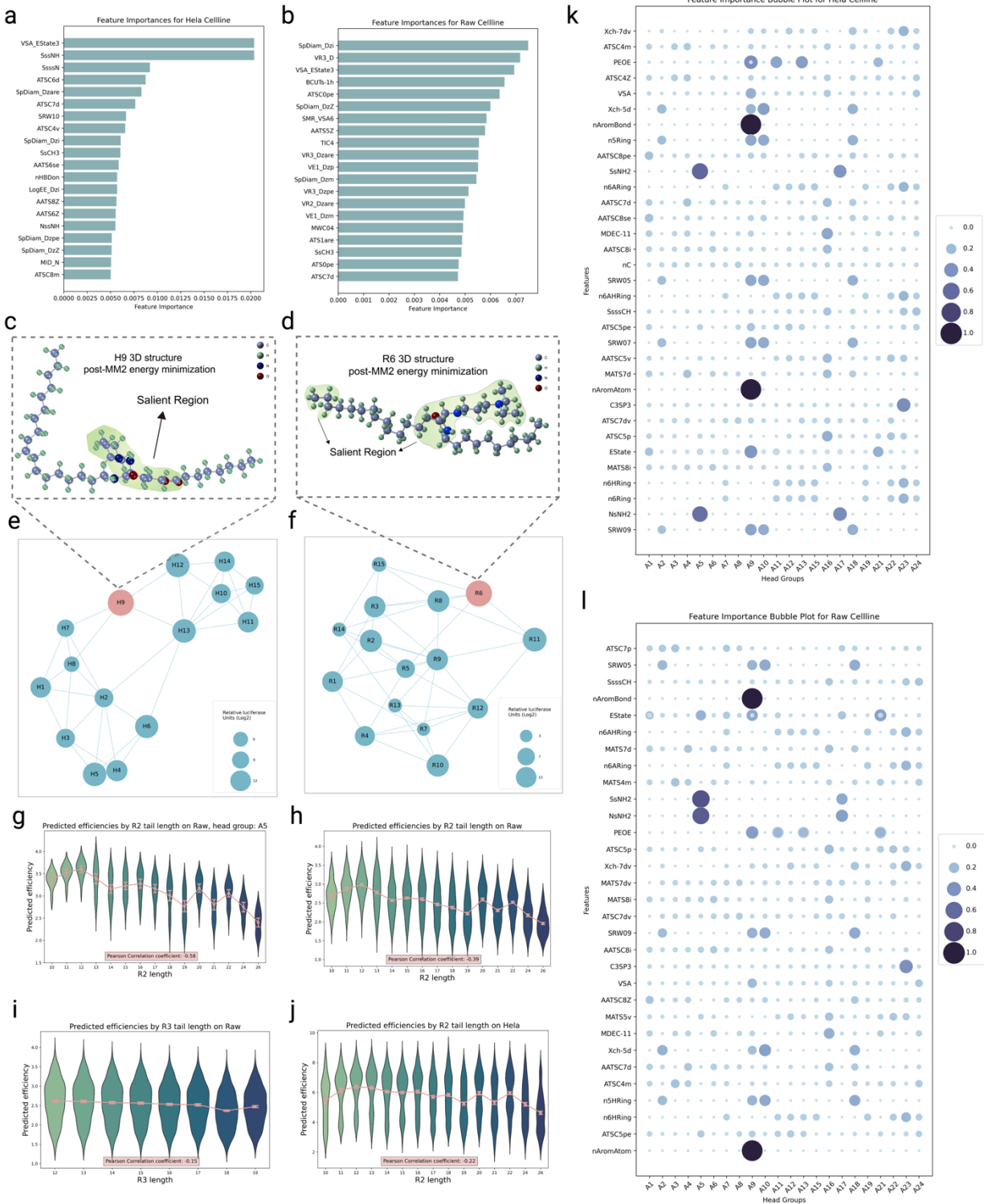




759

760 **Figure 4. Accelerating screening of new lipids for EGFP-mRNA delivery in macrophage**  
 761 **through the platform.** (a) The UMAP plot of the predicted molecule trans potencies. (b) Top  
 762 Head groups distribution and (c) top 3 and bottom 2 tail combinations distribution in RAW 264.7.  
 763 (d) Validate 15 lipid candidates for RAW 264.7 (n=2). (e) The top-performing formulation  
 764 parameters used in the optimization of R6 LNPs in RAW 264.7. (f) Comparison of the Fluc-mRNA  
 765 transfection potency of different LNPs in RAW 264.7 cells (n=6). (H9 LNPs, MC3 LNPs, R6  
 766 original screen formulation LNPs and optimized formulation LNPs). (g) Comparison of the  
 767 efficacy of LNPs (H9 LNPs and R6 LNPs) in HeLa cells (n=3). (h) Percentage of GFP positive  
 768 cells on RAW 264.7 after treatment with MC3 LNPs, H9 LNPs and H6 LNPs. Quantitative  
 769 analysis of flow cytometry data of RAW 264.7 cells (n=3). (i) Confocal images of RAW 264.7

770 cells transfected by GFP-mRNA LNPs. Green represents GFP, and blue represents the nucleus  
 771 (DAPI). Statistical significance was analyzed by the two-tailed Student's t-test. \* = p-value < 0.05,  
 772 \*\* = p-value < 0.01, \*\*\* = p-value < 0.005. Data are presented as mean ± SD.



774 **Figure 5. Model feature explanation and finding.** (a, b) Top 20 most important molecular  
775 descriptors identified by this model fine-tuned for HeLa and RAW 264.7 cell lines, respectively.  
776 (c, d) 3D visualization of H9 and R6 structures, respectively, with salient region highlighted. (e, f)  
777 Similarity networks for the 15 top lipid candidates in HeLa and RAW 264.7 cell lines respectively,  
778 with each candidate linked to its four closest neighbors. (g) Violin plot illustrating the distribution  
779 of predicted potencies across different R2 tail lengths, from LNPs of the top performing head group  
780 A5 for the RAW 264.7 cell line. (h) A similar violin plot as in (g), but focusing on LNPs of the  
781 entire candidate set. (i) A similar violin plot as in (h), but focusing on R3 tail lengths. (j) A similar  
782 violin plot as in (h), but focusing on LNPs of the entire candidate set for HeLa cell line. (k, l) Top  
783 2 most important molecular descriptors identified by this model fine-tuned for HeLa and RAW  
784 264.7 cell lines respectively, for each head group.

785

786

## Reference

- 787 1. Qin, S. et al. mRNA-based therapeutics: powerful and versatile tools to combat diseases. *Signal*  
788 *Transduction and Targeted Therapy* **7**, 166 (2022).
- 789 2. Hou, X., Zaks, T., Langer, R. & Dong, Y. Lipid nanoparticles for mRNA delivery. *Nature Reviews*  
790 *Materials* **6**, 1078-1094 (2021).
- 791 3. Kim, Y.-K. RNA therapy: rich history, various applications and unlimited future prospects.  
792 *Experimental & Molecular Medicine* **54**, 455-465 (2022).
- 793 4. Mendes, B.B. et al. Nanodelivery of nucleic acids. *Nature Reviews Methods Primers* **2**, 24 (2022).
- 794 5. Mitchell, M.J. et al. Engineering precision nanoparticles for drug delivery. *Nature Reviews Drug*  
795 *Discovery* **20**, 101-124 (2021).
- 796 6. Nasreen, S. et al. Effectiveness of COVID-19 vaccines against symptomatic SARS-CoV-2  
797 infection and severe outcomes with variants of concern in Ontario. *Nature microbiology* **7**, 379-  
798 385 (2022).
- 799 7. Patrignani, A. et al. Acute myocarditis following Comirnaty vaccination in a healthy man with  
800 previous SARS-CoV-2 infection. *Radiology Case Reports* **16**, 3321-3325 (2021).
- 801 8. Akinc, A. et al. The Onpattro story and the clinical translation of nanomedicines containing nucleic  
802 acid-based drugs. *Nature nanotechnology* **14**, 1084-1087 (2019).
- 803 9. Rüger, J., Ioannou, S., Castanotto, D. & Stein, C.A. Oligonucleotides to the (gene) rescue: FDA  
804 approvals 2017–2019. *Trends in pharmacological sciences* **41**, 27-41 (2020).
- 805 10. Chaudhary, N., Weissman, D. & Whitehead, K.A. mRNA vaccines for infectious diseases:  
806 principles, delivery and clinical translation. *Nature Reviews Drug Discovery* **20**, 817-838 (2021).
- 807 11. Kim, M. et al. Engineered ionizable lipid nanoparticles for targeted delivery of RNA therapeutics  
808 into different types of cells in the liver. *Science Advances* **7**, eabf4398 (2021).
- 809 12. Degors, I.M., Wang, C., Rehman, Z.U. & Zuhorn, I.S. Carriers break barriers in drug delivery:  
810 endocytosis and endosomal escape of gene delivery vectors. *Accounts of chemical research* **52**,  
811 1750-1760 (2019).
- 812 13. Wittrup, A. et al. Visualizing lipid-formulated siRNA release from endosomes and target gene  
813 knockdown. *Nature biotechnology* **33**, 870-876 (2015).
- 814 14. Xu, E., Saltzman, W.M. & Piotrowski-Daspi, A.S. Escaping the endosome: assessing cellular  
815 trafficking mechanisms of non-viral vehicles. *Journal of Controlled Release* **335**, 465-480 (2021).
- 816 15. Miao, L. et al. Delivery of mRNA vaccines with heterocyclic lipids increases anti-tumor efficacy  
817 by STING-mediated immune cell activation. *Nature biotechnology* **37**, 1174-1185 (2019).
- 818 16. Li, B. et al. Combinatorial design of nanoparticles for pulmonary mRNA delivery and genome  
819 editing. *Nature Biotechnology* (2023).
- 820 17. Han, X. et al. An ionizable lipid toolbox for RNA delivery. *Nat Commun* **12**, 7233 (2021).
- 821 18. Zador, A. et al. Catalyzing next-generation Artificial Intelligence through NeuroAI. *Nature*  
822 *Communications* **14**, 1597 (2023).
- 823 19. Bhardwaj, G. et al. Accurate de novo design of membrane-traversing macrocycles. *Cell* **185**, 3520-  
824 3532. e3526 (2022).
- 825 20. Yeh, A.H.-W. et al. De novo design of luciferases using deep learning. *Nature* **614**, 774-780 (2023).
- 826 21. Paul, D. et al. Artificial intelligence in drug discovery and development. *Drug Discov Today* **26**,  
827 80-93 (2021).
- 828 22. Melo, M.C.R., Maasch, J.R.M.A. & de la Fuente-Nunez, C. Accelerating antibiotic discovery  
829 through artificial intelligence. *Communications Biology* **4**, 1050 (2021).
- 830 23. Ma, Y. et al. Identification of antimicrobial peptides from the human gut microbiome using deep  
831 learning. *Nature Biotechnology* **40**, 921-931 (2022).
- 832 24. McCloskey, K. et al. Machine Learning on DNA-Encoded Libraries: A New Paradigm for Hit  
833 Finding. *Journal of Medicinal Chemistry* **63**, 8857-8866 (2020).
- 834 25. Stokes, J.M. et al. A Deep Learning Approach to Antibiotic Discovery. *Cell* **180**, 688-702.e613  
835 (2020).

- 836 26. Wang, W. et al. Prediction of lipid nanoparticles for mRNA vaccines by the machine learning  
837 algorithm. *Acta Pharmaceutica Sinica B* **12**, 2950-2962 (2022).
- 838 27. Huang, Y. et al. High-throughput microbial culturomics using automation and machine learning.  
839 *Nature Biotechnology*, 1-10 (2023).
- 840 28. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. in International conference on machine learning  
841 1597-1607 (PMLR, 2020).
- 842 29. Nazeri, M.T., Farhid, H., Mohammadian, R. & Shaabani, A. Cyclic Imines in Ugi and Ugi-Type  
843 Reactions. *ACS Combinatorial Science* **22**, 361-400 (2020).
- 844 30. Moriwaki, H., Tian, Y.-S., Kawashita, N. & Takagi, T. Mordred: a molecular descriptor calculator.  
845 *Journal of Cheminformatics* **10**, 4 (2018).
- 846 31. Yang, L. et al. Recent Advances in Lipid Nanoparticles for Delivery of mRNA. *Pharmaceutics* **14**  
847 (2022).
- 848 32. Barnard, J.M., Downs, G.M., von Scholley-Pfab, A. & Brown, R.D. Use of Markush structure  
849 analysis techniques for descriptor generation and clustering of large combinatorial libraries.  
850 *Journal of Molecular Graphics and Modelling* **18**, 452-463 (2000).
- 851 33. Kaczmarek, J.C. et al. Optimization of a degradable polymer-lipid nanoparticle for potent systemic  
852 delivery of mRNA to the lung endothelium and immune cells. *Nano letters* **18**, 6449-6454 (2018).
- 853 34. Eygeris, Y., Gupta, M., Kim, J. & Sahay, G. Chemistry of Lipid Nanoparticles for RNA Delivery.  
854 *Accounts of Chemical Research* **55**, 2-12 (2022).
- 855 35. Kim, M. et al. Engineered ionizable lipid nanoparticles for targeted delivery of RNA therapeutics  
856 into different types of cells in the liver. *Science Advances* **7**, eabf4398 (2021).
- 857 36. Miao, L. et al. Delivery of mRNA vaccines with heterocyclic lipids increases anti-tumor efficacy  
858 by STING-mediated immune cell activation. *Nature Biotechnology* **37**, 1174-1185 (2019).
- 859 37. Zhang, N.-N. et al. A Thermostable mRNA Vaccine against COVID-19. *Cell* **182**, 1271-  
860 1283.e1216 (2020).
- 861 38. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for  
862 dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- 863 39. Lam, K. et al. Unsaturated, Trialkyl Ionizable Lipids are Versatile Lipid-Nanoparticle Components  
864 for Therapeutic and Vaccine Applications. *Advanced Materials* **35**, 2209624 (2023).
- 865 40. Lee, S.M. et al. A systematic study of unsaturation in lipid nanoparticles leads to improved mRNA  
866 transfection in vivo. *Angewandte Chemie* **133**, 5912-5917 (2021).
- 867 41. Rhym, L.H., Manan, R.S., Koller, A., Stephanie, G. & Anderson, D.G. Peptide-encoding mRNA  
868 barcodes for the high-throughput in vivo screening of libraries of lipid nanoparticles for mRNA  
869 delivery. *Nature Biomedical Engineering* (2023).
- 870 42. Sedic, M. et al. Safety Evaluation of Lipid Nanoparticle-Formulated Modified mRNA in the  
871 Sprague-Dawley Rat and Cynomolgus Monkey. *Vet Pathol* **55**, 341-354 (2018).
- 872 43. Muzumdar, M.D., Tasic, B., Miyamichi, K., Li, L. & Luo, L. A global double-fluorescent Cre  
873 reporter mouse. *genesis* **45**, 593-605 (2007).
- 874 44. Boettler, T. et al. SARS-CoV-2 vaccination can elicit a CD8 T-cell dominant hepatitis. *J Hepatol*  
875 **77**, 653-659 (2022).
- 876 45. Seow, Y. & Wood, M.J. Biological Gene Delivery Vehicles: Beyond Viral Vectors. *Molecular*  
877 *Therapy* **17**, 767-777 (2009).
- 878 46. Kauffman, K.J. et al. Rapid, Single-Cell Analysis and Discovery of Vected mRNA Transfection  
879 In Vivo with a loxP-Flanked tdTomato Reporter Mouse. *Mol Ther Nucleic Acids* **10**, 55-63 (2018).
- 880 47. Kumar, A.R.K., Shou, Y., Chan, B., L., K. & Tay, A. Materials for Improving Immune Cell  
881 Transfection. *Advanced Materials* **33**, 2007421 (2021).
- 882 48. Van Hoeck, J., Braeckmans, K., De Smedt, S.C. & Raemdonck, K. Non-viral siRNA delivery to T  
883 cells: Challenges and opportunities in cancer immunotherapy. *Biomaterials* **286**, 121510 (2022).
- 884 49. Rampado, R. & Peer, D. Design of experiments in the optimization of nanoparticle-based drug  
885 delivery systems. *Journal of Controlled Release* **358**, 398-419 (2023).

- 886 50. Labute, P. A widely applicable set of descriptors. *Journal of Molecular Graphics and Modelling*  
887 **18**, 464-477 (2000).
- 888 51. Albertsen, C.H. et al. The role of lipid components in lipid nanoparticles for vaccines and gene  
889 therapy. *Advanced Drug Delivery Reviews*, 114416 (2022).
- 890 52. Boström, J., Brown, D.G., Young, R.J. & Keserü, G.M. Expanding the medicinal chemistry  
891 synthetic toolbox. *Nature Reviews Drug Discovery* **17**, 709-727 (2018).
- 892 53. Zhang, M. et al. A survey on graph diffusion models: Generative ai in science for molecule, protein  
893 and material. *arXiv preprint arXiv:2304.01565* (2023).
- 894 54. Hoogeboom, E., Satorras, V.G., Vignac, C. & Welling, M. in International Conference on Machine  
895 Learning 8867-8887 (PMLR, 2022).
- 896 55. Landrum, G. Rdkit: Open-source cheminformatics software. (2016).
- 897 56. Xu, K., Hu, W., Leskovec, J. & Jegelka, S. How powerful are graph neural networks? *arXiv*  
898 *preprint arXiv:1810.00826* (2018).
- 899 57. Wang, Y., Wang, J., Cao, Z. & Barati Farimani, A. Molecular contrastive learning of  
900 representations via graph neural networks. *Nature Machine Intelligence* **4**, 279-287 (2022).
- 901 58. Kingma, D.P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*  
902 *arXiv:1412.6980* (2014).
- 903 59. Sundararajan, M., Taly, A. & Yan, Q. in International conference on machine learning 3319-3328  
904 (PMLR, 2017).
- 905 60. Kokhlikyan, N. et al. Captum: A unified and generic model interpretability library for pytorch.  
906 *arXiv preprint arXiv:2009.07896* (2020).
- 907 61. Wellawatte, G.P., Seshadri, A. & White, A.D. Model agnostic generation of counterfactual  
908 explanations for molecules. *Chemical science* **13**, 3697-3705 (2022).

909