

Atom-by-atom protein generation and beyond with language models

Daniel Flam-Shepherd,^{1,2} Kevin Zhu,^{1,2} and Alán Aspuru-Guzik^{1,2,3,4,5}

¹*Department of Computer Science, University of Toronto, Toronto, Ontario M5S 2E4, Canada*

²*Vector Institute for Artificial Intelligence, Toronto, Ontario M5S 1M1, Canada*

³*Department of Chemistry, University of Toronto, Toronto, Ontario M5G 1Z8, Canada*

⁴*Canadian Institute for Advanced Research, Toronto, Ontario M5G 1Z8, Canada*

⁵*Acceleration Consortium, Toronto, Ontario M5S 3H6, Canada*

Protein language models learn powerful representations directly from sequences of amino acids. However, they are constrained to generate proteins with only the set of amino acids represented in their vocabulary. In contrast, chemical language models learn atom-level representations of smaller molecules that include every atom, bond, and ring. In this work, we show that chemical language models can learn atom-level representations of proteins enabling protein generation unconstrained to the standard genetic code and far beyond it. In doing so, we show that language models can generate entire proteins atom by atom—effectively learning the multiple hierarchical layers of molecular information that define proteins from their primary sequence to their secondary, and tertiary structure. We demonstrate language models are able to explore beyond protein space—generating proteins with modified sidechains that form unnatural amino acids. Even further, we find that language models can explore chemical space and protein space simultaneously and generate novel examples of protein-drug conjugates. The results demonstrate the potential for biomolecular design at the atom level using language models.

Proteins are essential components of all life on Earth and are involved in every cellular process. As a result, protein engineering is one of the most important areas of scientific discovery. Significant progress has been made and proteins have been engineered for therapies against viruses [1] and cancer [2], as well as to alleviate genetic diseases directly [3–5]. Artificial intelligence has enormous potential to accelerate scientific progress and automate protein engineering. Already it has led to a breakthrough in highly accurate protein structure prediction [6]. In particular, language models have already begun to have a major impact on protein design [7–9].

The important functions proteins carry out and the structure responsible for them originate in the patterns of amino acids in their primary sequence. Indeed, most language models represent proteins using sequences of amino acids [7–9]. However, this ignores atom-level interactions, precluding the model from representing any atom-level protein modification. Allowing for atom-level representations would enable protein generation outside of the genetic code and allow language models to explore an expanded space of biomolecules. Specifically, this would make it possible for the model to propose new unnatural side chains, attach small molecules, and generate linkers between residues that form large macrocycles.

To learn atom-level representations we must turn to chemical language models, however, these models are typically used for smaller drug-like molecules. Similar to their protein variants—chemical language models are deep neural networks trained using masking or next-token prediction [10] but use atom-level linear sequences parsed from molecular graphs [10, 11]. These sequences completely represent the molecule including all atoms, bonds, rings, aromaticity, branching, and stereochemistry. The two most prominent sequence representations are SMILES strings [12] or SELFIES strings [13] which

are completely robust and always valid.

Recently, chemical language models [10] were found to have the ability to generate larger, complex molecules, relative to small drug-like molecules such as the largest molecules in PubChem. These molecules are still much smaller than proteins, but this indicates that atom-level protein generation with language models is feasible. In this work, we demonstrate that chemical language models are capable of generating entire proteins atom by atom, including biomolecules beyond protein space. Specifically, we train models on various biomolecules including proteins from the protein databank. We also create two other synthetic biomolecular datasets, first modified proteins with unnatural amino acids, and proteins with small molecule attachments—specifically single domain antibodies (sdAbs) from the antibody structural database [14] attached to molecules from the ZINC dataset [15].

We discover that chemical language models can learn the language of proteins entirely from scratch—by learning to generate atom-level sequences that define proteins with valid primary sequences that correspond to meaningful secondary, and tertiary structure, which we check using AlphaFold [6] structure predictions. Importantly, the language model learns valid protein backbones and natural amino acid structures as well as the primary sequence patterns in the training proteins. We further demonstrate that language models can generate beyond the standard genetic code—proteins with novel sidechains that are more complex than the set of standard amino acids. Additionally, we also show that chemical language models can generate novel proteins and small molecules together at the same time as protein-drug conjugates. In particular, we find that the model learns both the protein space of the single domain antibodies and the chemical space defined by the ZINC molecules—generating antibody-drug conjugates with valid and novel pro-

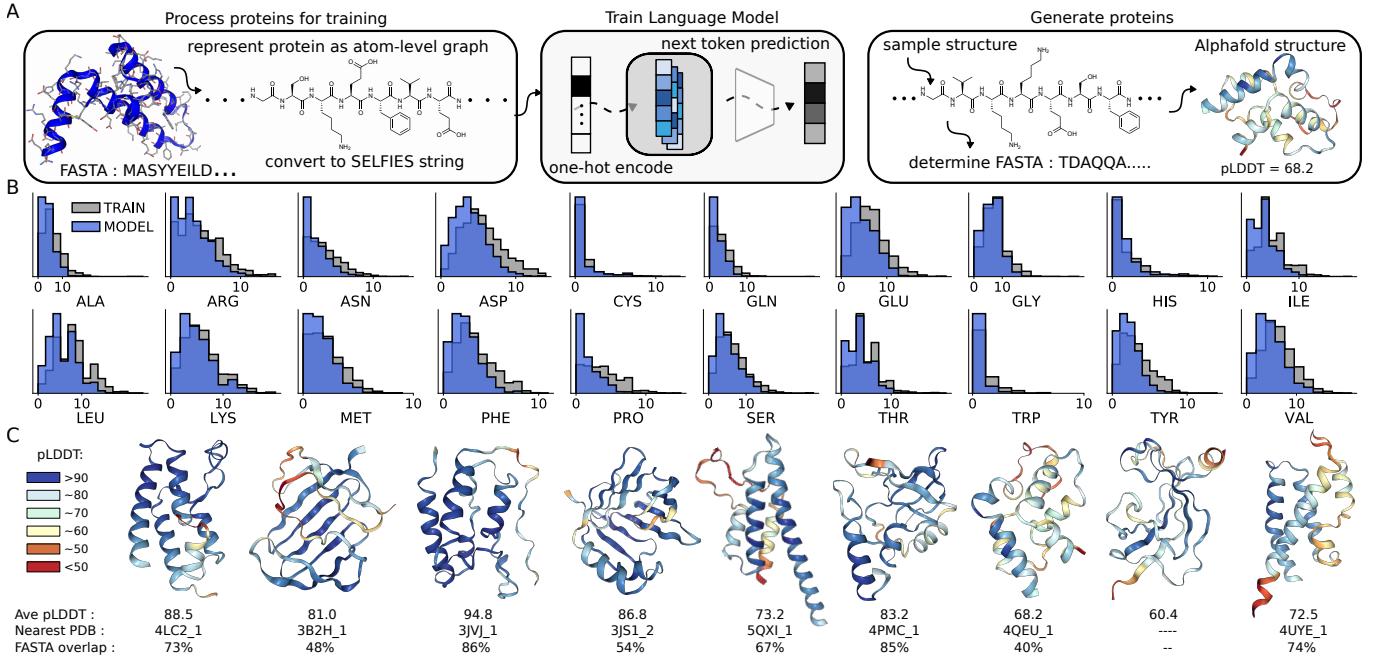


FIG. 1. **Proteins** (A) Dataset preparation. The training workflow for the model: training, generation, amino acid sequence determination, and AlphaFold visualization. (B) Comparison of amino acid distributions. (C) AlphaFold visualizations of model-generated proteins coloured by pLDDT, including the PDB ID of the closest protein and its % sequence overlap.

tein sequences and structures attached to novel drug-like molecules warheads similar to the structures in ZINC.

RESULTS

In this study, the datasets are constructed by using small proteins from the Protein Data Bank (PDB), specifically between 50 and 150 residues. We use atom-level graph representations of each protein so that sidechain modifications can be made directly. For training, each protein can be parsed to a linear string representation, and random data augmentation can be used to increase the training data size. We describe the main details and results for each dataset in the following sections.

A. Proteins

For the first dataset, which consists of proteins with standard amino acids and no sidechain modifications, we test the ability of the language model to explore protein space while maintaining protein structure and constraints. After training, as shown in Figure 1(A))– we generate a thousand (1K) samples from the language model and evaluate their atom, residue, and protein-level properties. At the protein level, we determine if generated samples are proteins and attempt to determine their primary sequence. If we can ascertain their primary sequence, we can use AlphaFold2 [6] to further evaluate

if the model has learned the amino acid sequences that correspond to good structure predictions. Additionally, we study model samples for their distribution of amino acids and other atom-level properties that can be computed using rdkit [16].

We check if molecules generated by the model are actually proteins by analyzing if they preserve the basic structure of the protein backbone and natural amino acids form. First, we perform a backbone structure search and then attempt to arrange the backbone from the N terminus to the C terminus while simultaneously classifying each sidechain using another substructure search for the standard set of amino acids. If this is successful and there are no discontinuities in the backbone or other side chain errors, then we classify the sample as a protein and parse the amino acid sequence. By this process, we determine roughly $\sim 68.2\%$ of samples are proteins, furthermore, all the parsed amino acid sequences are unique (there are no duplicates and the model isn't repeating specific proteins) and novel (they are different from the training sequences).

We compare the distribution of amino acids in the training sequences to the distribution learned by the model based on the generated samples. We plot histograms, in Fig. 1(B), displaying the frequency of occurrence of every amino acid in samples from both the model and the training data– both distributions are very similar and mostly overlap but for some amino acids, the language model slightly underestimates the training frequencies.

Using Alphafold [6], in Fig. 1(B), we visualize se-

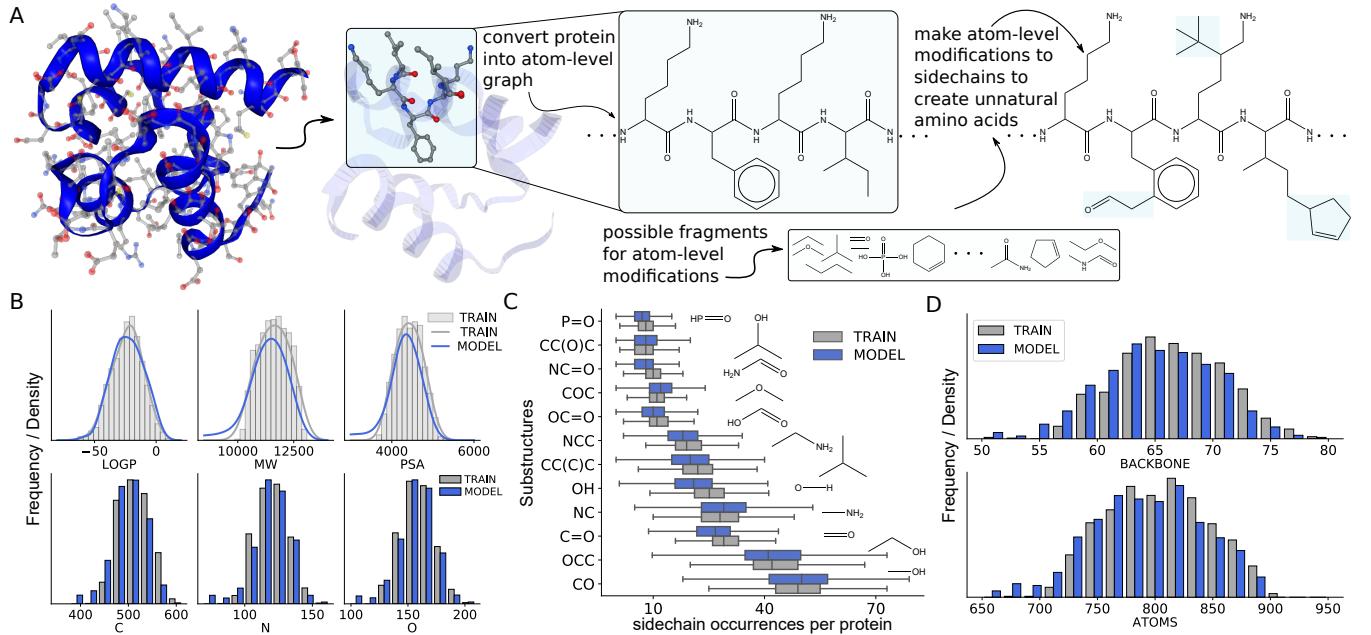


FIG. 2. Proteins with unnatural amino acids (A) Dataset of unnatural proteins built by random side chain modification. (B) Histograms and density plots of atom-level measures or properties. Density plots use a Gaussian kernel density estimator (KDE) that is fit to LogP, MW, and PSA values of the training molecules by tuning the bandwidth parameter. (C) Boxplots measuring the number of occurrences of small fragments in sidechains of model and training proteins. (D) Histograms of backbone size and heavy atom number per protein for model and training proteins.

lected examples of proteins generated by the language model. In each sample, residues are color-coded according to pLDDT, which is a per-residue estimate of the model's confidence on a scale from 0 to 100. Regions with pLDDT > 90 are dark blue and have high accuracy. Regions with pLDDT from 90 down to 70 are still expected to be good predictions and are colored light blue that transitions to green (with decreasing confidence). Regions with pLDDT between 50 and 70 are lower confidence and are colored yellow to green. The regions with pLDDT < 50 are not confident and likely disordered—these are colored red.

On this scale, in Fig. 1(C), we see that the proteins that are generated by the model result in good structure predictions—ranging between 70 and 90 pLDDT. This indicates that the model can generate proteins with well-defined structures that are not disordered. For a simple baseline comparison, we considered sequences of random amino acids, the structure predictions for these consistently result in disordered proteins with low pLDDT < 50.

Additionally, in Fig. 1(C), the proteins generated by the language model contain a variety of secondary structures including alpha helices, beta sheets, and omega loops. Globally, the generated proteins combine many of these secondary structures into various and unique domains. We can conclude, based on these samples, and further examples in Supplementary Fig. S1, that language models can generate proteins, atom by atom, not just with valid primary sequences but proteins with

meaningful secondary and tertiary structure.

Furthermore, the generated proteins are similar to their nearest training examples in the PDB, to show this in Fig. 1(C) and Supplementary Fig. S1, under each protein we label the primary sequence percentage overlap between the generated proteins and their most similar PDB training example—which ranges from 86% to 40% with one other generated protein that has no nearest PDB training example. Based on this, it is evident that the model draws heavily from the amino acid sequence patterns in its training data but does not memorize them.

Also, in Supplementary Fig. S2, we plot histograms comparing atom-level properties of the samples generated from the model with the training data. The model roughly approximates the training distribution of atoms but slightly underestimates some properties.

B. Proteins with unnatural amino acids

The next dataset, whose construction is depicted in Fig 2(A), consists of proteins that have random sidechain modifications creating proteins with unnatural amino acids. For this dataset, we select a subset of smaller proteins with 50 to 80 residues from the previous protein dataset. We then modify the selected protein by attaching a randomly chosen small fragment using a random attachment point on every sidechain. This produces a dataset of proteins that are entirely comprised of "unnatural" amino acids. We train language models on these

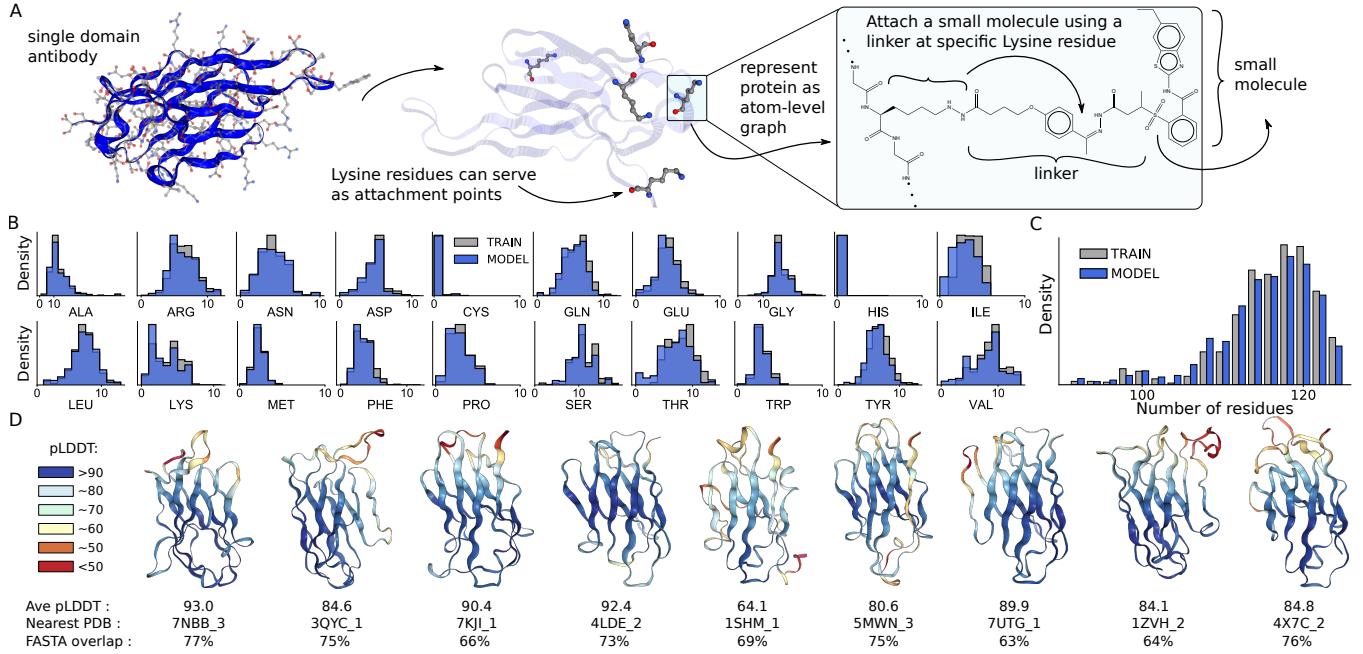


FIG. 3. Antibody Drug Conjugates – sdAbs (A) Single domain Antibody-drug conjugates dataset creation overview. (B) Comparison of amino acid distributions for training and model sdAbs. (C) Histogram comparing the size of training and model sdAbs (by number of residues). (D) Example sdAbs (with warheads excluded) generated by the language model visualized by Alphafold and coloured by pLDDT. Under each, we include the PDB ID of the closest protein and its % sequence overlap.

unnatural proteins in order to test the ability of chemical language models to generate biomolecules beyond protein space. Further details about building the training data can be found in the Methods section. Additionally, an example of the unnatural amino acids from a single training protein and model protein can be found in Supplementary Fig. S3.

After training, we again generate 1K samples from the language model for evaluation– the results are shown in Fig 2, where we test the model’s ability to capture atom-level and sidechain-level properties of the unnatural proteins. First, in Fig 2 (B), we see that the model learns the continuous atom-level properties of the training proteins including octanol-water partition coefficient (LogP) [17], exact molecular weight (MW) and the topological polar surface area (PSA), in addition to learning the number of carbon, nitrogen, and oxygen. Then in Fig 2 (C), we see the model learns a similar sidechain structure to the training sidechains as determined by a structure search in each sidechain for a basic set of small fragments. Lastly, in Fig 2 (D), the model learns to generate unnatural proteins of similar atom number and backbone size to the training proteins– but does tend to slightly underestimate the size of both.

C. Antibody Drug Conjugates

Next, we test the ability of the language model to generate proteins attached to small molecules and simulta-

neously explore protein space and chemical space. One of the most promising examples of this kind of biomolecule with immense therapeutic potential are antibody-drug conjugates, which are a form of cancer therapy intended to target and kill cancer cells but spare healthy cells [18]. Structurally, they are composed of an antibody attached to single or multiple anticancer drugs typically using some linker molecule. To construct a synthetic dataset of antibody-drug conjugates, as shown in Fig 3 (A), we attach a single drug-like molecule from the ZINC dataset [15] to single-domain antibodies (sdAbs) from the structural antibody dataset [14, 19] in order to test the ability of language models to generate antibody-drug conjugates. We use two possible linkers for cysteine attachments and two other linkers for lysine attachments. Both linkers are selected from real antibody-drug conjugates described in [18]. The linker is randomly attached to the small molecule from ZINC and the specific lysine or cysteine residue for attachment is also randomly chosen. Since there are only 1K sdAbs in the structural antibody dataset we use data augmentation to expand the dataset size to 250K proteins that can be attached to every molecule in ZINC.

After training, we again generate 1K samples from the language model for evaluation, We first test the model’s ability to explore protein space and learn the distribution of single-domain antibodies– the results are shown in Fig 3 (B-D). Similar to the standard protein data, we compare the distribution of amino acids in the training sequences to the distribution learned by the model. We

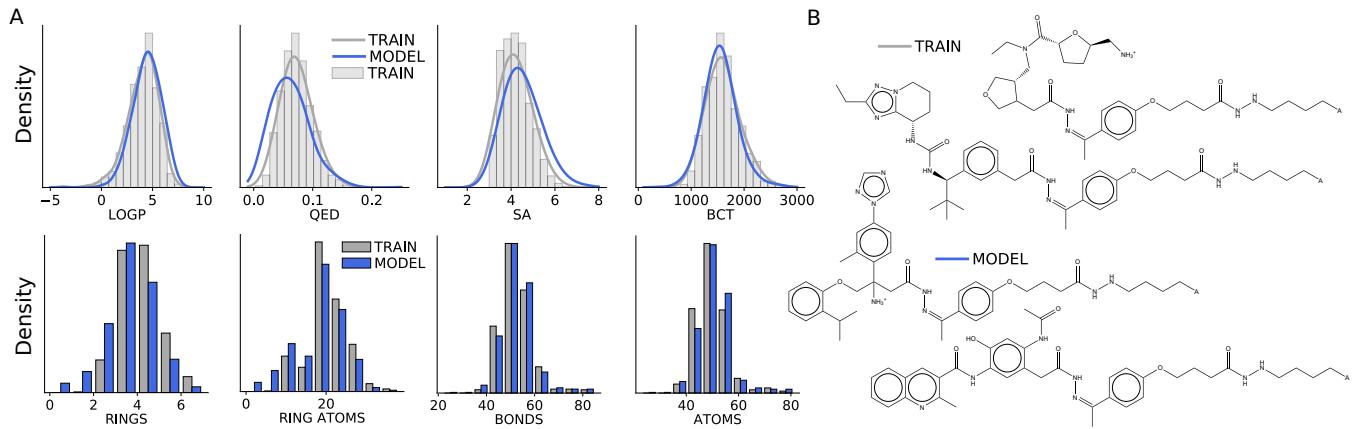


FIG. 4. Antibody Drug Conjugates – Warheads (A) Histograms and density plots of atom-level measures or properties of warheads. Density plots use a Gaussian kernel density estimator fit to LogP, MW, and PSA values for the training molecules by tuning the bandwidth parameter. (B) Examples of model and train “warheads”.

plot histograms, in Fig. 4(B), displaying the frequency of occurrence of every amino acid in samples from both the model and the training data– from these, we can see the language model accurately learns the training distribution of amino acids. Similarly, in Fig. 4(C), the model accurately learns the size of the training sdAbs.

Similar to the standard proteins, we can attempt to determine the amino acid sequences of the single-domain antibodies (ignoring the warheads). We determine roughly $\sim 90.8\%$ of samples are proteins and their primary sequences are entirely unique and novel (there are no duplicates and all are different from training sequences).

Even further, AlphaFold structure predictions [6], visualized in Fig. 3(D) and Supplementary Fig. S7, confidently show that the language model produces sequences that fold into the expected structure for single domain antibodies. Additionally, based on the primary sequence overlap of model samples with their nearest PDB training example in Fig. 1(C) and Supplementary Fig. S1, the model, without memorizing, learns the amino acid sequence structure that defines the training sdAbs. The primary sequence overlap with training examples ranges from 63% to 93% in the supplementary information. Investigating further, we see that the model draws heavily from the sdAb sequences making new examples of sequences by memorizing small snippets of amino acids and using larger training snippets but with a large number of single mutations randomly distributed throughout the snippet.

From the training examples and model samples, we detach and collect “warheads” which we expand the definition of to include the linker and sidechain in addition to the small molecule (warhead typically refers to just the small molecule). In Fig. 4 (B), two examples of train and model warheads are shown as graphs to clarify this. Additional model and training warheads are shown as graphs in Supplementary Fig. S9 and S8– as expected the same four linkers repeat across samples but the small

molecules attached to them differ and are structurally similar to the ZINC molecules in the training warheads.

We also evaluate the language model’s warheads in terms of their atom-level properties. In Fig 4 (A), the model captures the atom-level properties of the training warheads, specifically, it learns the continuous atom-level properties of the training warheads including LogP [17], drug-likeness (QED) [20], Synthetic Accessibility Score (SA) and molecular graph complexity (BCT) as well as the number of atoms, bonds, rings and atoms in rings. However, the model slightly underestimates the main modes for QED and SA as well as the number of rings per warhead.

Additionally, we assess the model warheads and compare them with the training warheads, we find that model warheads are unique (there are no duplicates and the model is not repeating a few examples) as well as novel (the model does not make exact copies of warheads from the training data). Given that the linkers are memorized, this indicates that the model is learning to generate new small molecules similar to ZINC molecules and effectively exploring chemical space at the same time it learns to explore the protein space defined by the sdAbs.

Also, in Supplementary Fig. S2, we see that the model does learn the atom-level properties of the training antibody drug-conjugates. Additionally, in Supplementary Fig. S4-S6, we show a single train antibody drug-conjugate and four model samples.

DISCUSSION

In this work, we show that chemical language models can generate, atom by atom, entire proteins, unnatural proteins, and protein drug conjugates. By analyzing generated samples we find that language models learn multiple hierarchical layers of molecular information that define the training biomolecules. This includes atom-level molecular properties or residue-level constraints for

backbone and amino acid structure as well as primary sequence patterns and motifs that define meaningful secondary and tertiary structure. Indeed, chemical language models learn to generate protein structures as sequence representations of atom-level graphs that are similar to the training proteins in the PDB.

Effectively we demonstrate that chemical language models can also serve as biological language models—capable of learning the language of proteins atom by atom. Importantly, in contrast to protein language models that only learn representations of amino acid sequences, chemical language models generate entire molecular graphs, and because of this we are able to show that language models can be used to explore not just chemical space and in between chemical and protein space but also protein space itself, beyond protein space, or even both chemical and protein space at the same time.

Further work should be done to ensure the model generates valid protein structure including correct backbone and amino acid form. This will also assist the model in learning distributions consisting of larger biomolecules including structures with more than 150 residues and multiple domains. Using memorizing Transformers [21] may help the model generate valid protein sequences. Also, other architectures built for longer sequence lengths [22] can increase the size and range of structures that the model can learn. Another limitation is that we do not consider the three-dimensional structure of the biomolecules and generate atom-level sequence representations. This problem can not be easily rectified because no training data with 3D information exists for unnatural proteins and protein-drug conjugates. A potential solution that does not require training data would be to use reinforcement learning [23] or bayesian optimization [24] and guide the model to generate 3D structure using energy.

The goal of this work is to demonstrate the power of chemical language models and their ability to learn atom-level representations of biomolecules. We envision future language models will be able to explore any combinatorial space in chemistry or biology using any representation type the user wishes [25].

I. METHODS

A. Datasets

From the PDB we successfully parse around $\sim 10K$ proteins between 50 and 150 residues. In all datasets, we only parse proteins that conform to atom-level graphs with no more than 2 macrocycles (created by residue-residue connections) – this makes primary sequence determination more successful. Given this constraint, we parse around $\sim 10K$ and $\sim 5K$ proteins from the PDB for the first two training datasets. In order to increase the size of the training data, we randomize the atom orderings of each protein in `rdkit` to obtain multiple different

random copies of each biomolecule as SMILES (and then SELFIES strings). Using this data augmentation we expand all training datasets to around $\sim 250K$ sequences. We use `rdkit` [16] to represent each protein as atom-level graphs and make side-chain modifications. We also made use of Colabfold [26] for quick visualization and NGLview [27] for figure construction.

B. Tokenization

We use SELFIES [13] version one for the sequence representation of atom-level protein graphs. Other than special tokens like `[BOS]`, `[EOS]`, `[PAD]`, `[UNK]`, the vocab \mathcal{T} consists of standard selfies tokens, encoding all information in a molecular graph including: atom tokens $\{[C], [N], \dots\}$, bond tokens $\{[=C], [#N], \dots\}$, ring tokens: $\{[Ring1], [Ring2], \dots\}$ branching tokens: $\{[Branch1_1], [Branch1_2], \dots\}$. In total for all datasets, the vocabulary is around ~ 30 tokens.

C. Language Modeling for Molecular Design

In language modeling for molecular design, we want to estimate the unsupervised distribution of the training molecules $(MOL_1, MOL_2, \dots, MOL_n)$ each composed of variable length sequences of tokens from a chemical language $[CT]_i$ where $CT \in \mathcal{T}$ such that $MOL = ([CT]_1, [CT]_2, \dots, [CT]_n)$. The joint probabilities over a single molecule can be written as

$$p(MOL) = \prod_{i=1}^n p([CT]_n | [CT]_{n-1}, \dots, [CT]_1) \quad (1)$$

These probabilities $p([CT]_n | [CT]_{n-1}, \dots, [CT]_1)$ are modeled using a Transformer [28] that is trained using stochastic gradient descent.

D. Training

During training, we one-hot encode SELFIES sequences using a basic vocabulary that consists of 30 possible alphabet tokens. All language models are trained using next-token prediction conditioned on the entire sequence for context. The training data only uses sequences that have a maximum length of 1664 tokens. We trained language models with decoder only, GPT-like architecture [29] with 4 attention heads and between 1 and 10 Million parameters. Language models are implemented in Python 3 with PyTorch [30]. Molecules properties are computed using `rdkit` [16].

II. ACKNOWLEDGEMENTS

A.A.-G. acknowledge funding from Dr. Anders G. Frøseth. A.A.-G. also acknowledges support from the Canada 150 Research Chairs Program, the Canada Industrial Research Chair Program, and from Google, Inc.

-
- Models were trained using the Canada Computing Systems [31]. This research was undertaken thanks in part to funding provided to the University of Toronto's Acceleration Consortium from the Canada First Research Excellence Fund.
- [1] R. Diskin, J. F. Scheid, P. M. Marcovecchio, A. P. West Jr, F. Klein, H. Gao, P. N. Gnanapragasam, A. Abadir, M. S. Seaman, M. C. Nussenzweig, et al., "Increasing the potency and breadth of an hiv antibody by using structure-based rational design," *Science* **334**, 1289 (2011).
- [2] G. P. Adams and L. M. Weiner, "Monoclonal antibody therapy of cancer," *Nature biotechnology* **23**, 1147 (2005).
- [3] A. V. Anzalone, P. B. Randolph, J. R. Davis, A. A. Sousa, L. W. Koblan, J. M. Levy, P. J. Chen, C. Wilson, G. A. Newby, A. Raguram, et al., "Search-and-replace genome editing without double-strand breaks or donor dna," *Nature* **576**, 149 (2019).
- [4] N. M. Gaudelli, A. C. Komor, H. A. Rees, M. S. Packer, A. H. Badran, D. I. Bryson, and D. R. Liu, "Programmable base editing of a•t to g•c in genomic dna without dna cleavage," *Nature* **551**, 464 (2017).
- [5] A. C. Komor, Y. B. Kim, M. S. Packer, J. A. Zuris, and D. R. Liu, "Programmable editing of a target base in genomic dna without double-stranded dna cleavage," *Nature* **533**, 420 (2016).
- [6] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, et al., "Highly accurate protein structure prediction with alphafold," *Nature* **596**, 583 (2021).
- [7] E. C. Alley, G. Khimulya, S. Biswas, M. AlQuraishi, and G. M. Church, "Unified rational protein engineering with sequence-based deep representation learning," *Nature methods* **16**, 1315 (2019).
- [8] B. L. Hie, V. R. Shanker, D. Xu, T. U. Bruun, P. A. Weidenbacher, S. Tang, W. Wu, J. E. Pak, and P. S. Kim, "Efficient evolution of human antibodies from general protein language models," *Nature Biotechnology* (2023).
- [9] A. Madani, B. Krause, E. R. Greene, S. Subramanian, B. P. Mohr, J. M. Holton, J. L. Olmos Jr, C. Xiong, Z. Z. Sun, R. Socher, et al., "Large language models generate functional protein sequences across diverse families," *Nature Biotechnology* , 1 (2023).
- [10] D. Flam-Shepherd, K. Zhu, and A. Aspuru-Guzik, "Language models can learn complex molecular distributions," *Nature Communications* **13**, 3293 (2022).
- [11] R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams, and A. Aspuru-Guzik, "Automatic chemical design using a data-driven continuous representation of molecules," *ACS central science* **4**, 268 (2018).
- [12] D. Weininger, "Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules," *Journal of chemical information and computer sciences* **28**, 31 (1988).
- [13] M. Krenn, F. Häse, A. Nigam, P. Friederich, and A. Aspuru-Guzik, "Selfies: a robust representation of semantically constrained graphs with an example application in chemistry," *arXiv preprint arXiv:1905.13741* (2019).
- [14] C. Schneider, M. I. J. Raybould, and C. M. Deane, "SAbDab in the age of biotherapeutics: updates including SAbDab-nano, the nanobody structure tracker," *Nucleic Acids Research* **50**, D1368 (2021), <https://academic.oup.com/nar/article-pdf/50/D1/D1368/42058451/gkab1050.pdf>.
- [15] J. J. Irwin and B. K. Shoichet, "Zinc- a free database of commercially available compounds for virtual screening," *Journal of chemical information and modeling* **45**, 177 (2005).
- [16] G. Landrum, "Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling," (2013).
- [17] S. A. Wildman and G. M. Crippen, "Prediction of physicochemical parameters by atomic contributions," *Journal of chemical information and computer sciences* **39**, 868 (1999).
- [18] A. Beck, L. Goetsch, C. Dumontet, and N. Corvaia, "Strategies and challenges for the next generation of antibody-drug conjugates," *Nature reviews Drug discovery* **16**, 315 (2017).
- [19] C. Schneider, M. I. Raybould, and C. M. Deane, "Sabdab in the age of biotherapeutics: updates including sabdab-nano, the nanobody structure tracker," *Nucleic acids research* **50**, D1368 (2022).
- [20] G. R. Bickerton, G. V. Paolini, J. Besnard, S. Muresan, and A. L. Hopkins, "Quantifying the chemical beauty of drugs," *Nature chemistry* **4**, 90 (2012).
- [21] Y. Wu, M. N. Rabe, D. Hutchins, and C. Szegedy, in *International Conference on Learning Representations* (2021).
- [22] R. Child, S. Gray, A. Radford, and I. Sutskever, "Generating long sequences with sparse transformers," *arXiv preprint arXiv:1904.10509* (2019).
- [23] D. Flam-Shepherd, A. Zhigalin, and A. Aspuru-Guzik, "Scalable fragment-based 3d molecular design with reinforcement learning," *arXiv preprint arXiv:2202.00658* (2022).
- [24] T. C. Wu, D. Flam-Shepherd, and A. Aspuru-Guzik, "Bayesian variational optimization for combinatorial spaces," *arXiv preprint arXiv:2011.02004* (2020).
- [25] D. Flam-Shepherd and A. Aspuru-Guzik, "Language models can generate molecules, materials, and protein binding sites directly in three dimensions as xyz, cif, and pdb files," *arXiv preprint arXiv:2305.05708* (2023).
- [26] M. Mirdita, K. Schütze, Y. Moriwaki, L. Heo, S. Ovchinnikov, and M. Steinegger, "Colabfold: making protein folding accessible to all," *Nature methods* **19**, 679 (2022).

- [27] H. Nguyen, D. A. Case, and A. S. Rose, “Nglview—interactive molecular graphics for jupyter notebooks,” *Bioinformatics* **34**, 1241 (2018).
- [28] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving language understanding by generative pre-training,” () .
- [29] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., “Language models are unsupervised multitask learners,” () .
- [30] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., “Pytorch: An imperative style, high-performance deep learning library,” *Advances in neural information processing systems* **32** (2019).
- [31] S. Baldwin, in *Journal of Physics: Conference Series*, Vol. 341 (IOP Publishing, 2012) p. 012001.

SUPPLEMENTARY INFORMATION

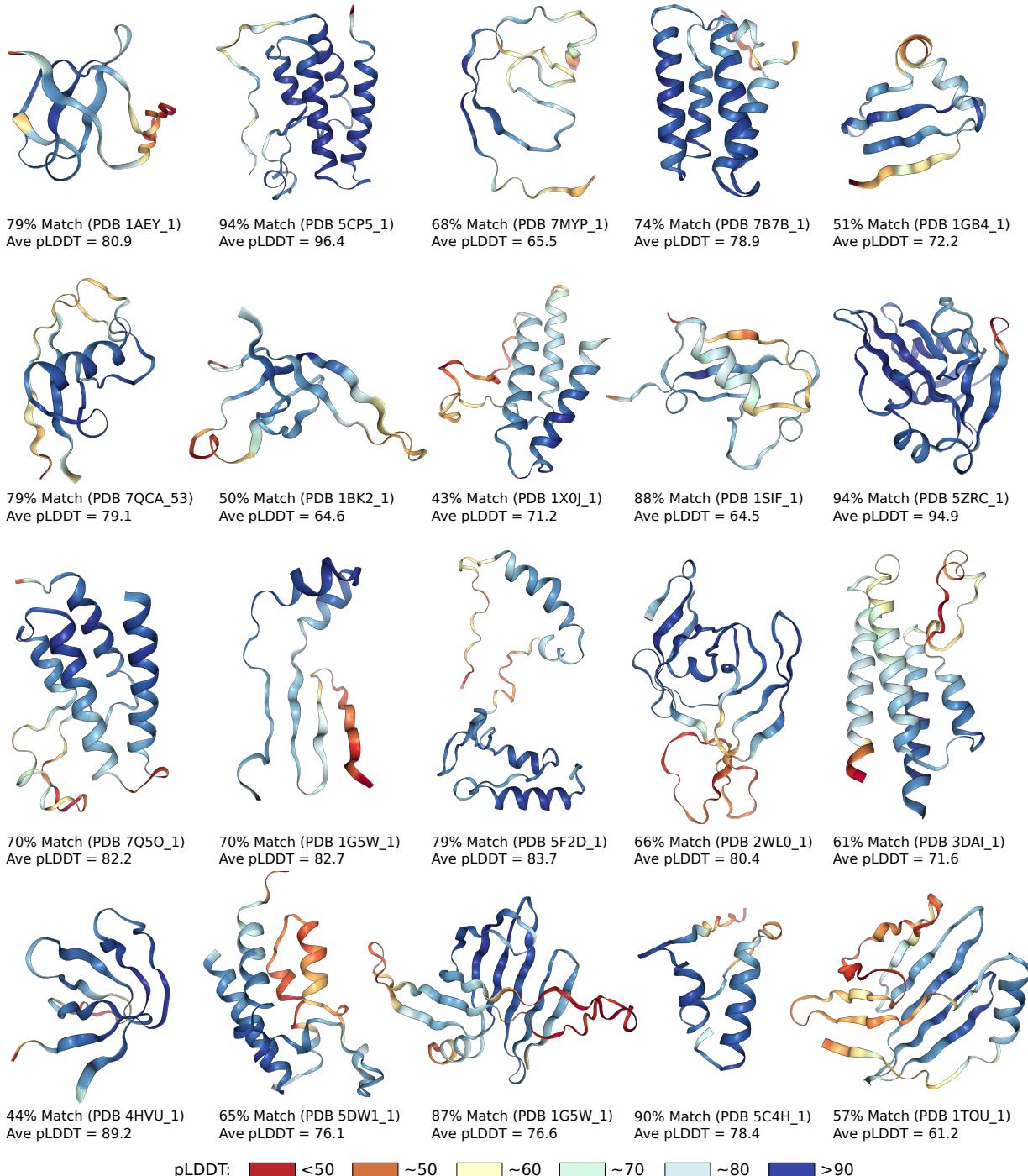


FIG. S1. Examples of proteins generated by the model visualized by AlphaFold and colored by pLDDT.

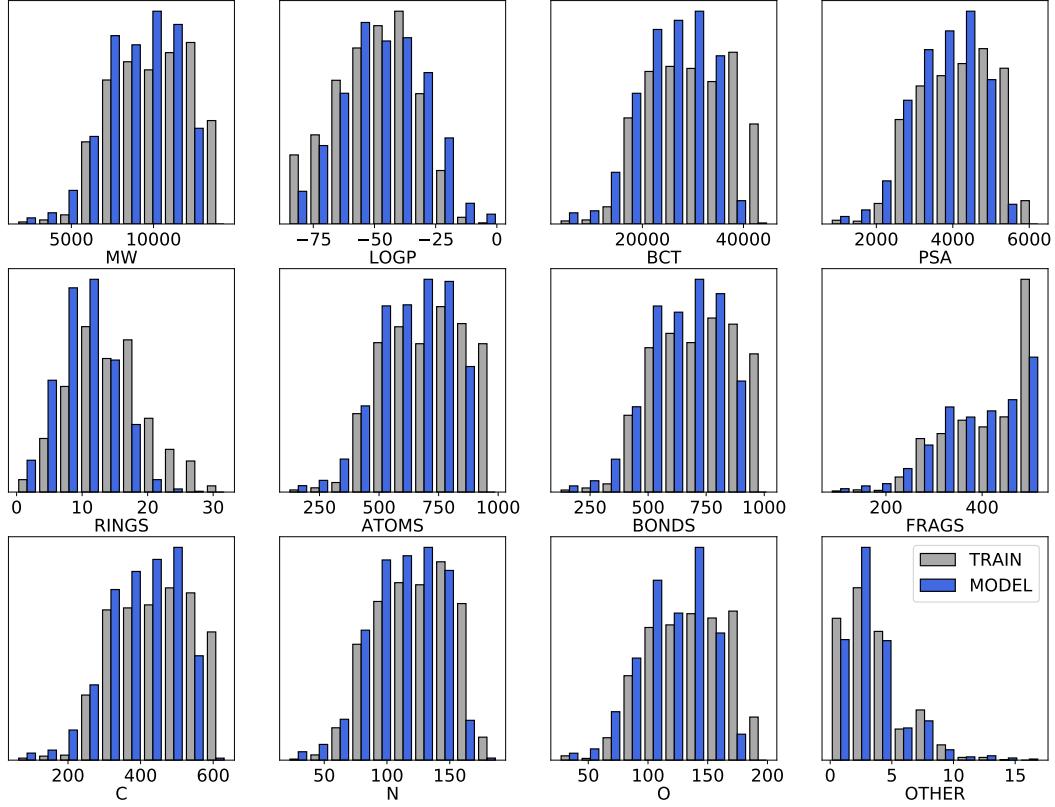


FIG. S2. Histograms of atom level properties for the basic proteins training data. These include exact molecular weight (MW), octanol-water partition coefficient (LogP) [17], molecular complexity (BCT), topological polar surface area (PSA), number of rings (RINGS), number of atoms (ATOMS), number of bonds (BONDS), number of fragments found by breaking up the molecule at rotatable bonds (FRAGS), number of carbons (C), number of nitrogens (N), number of oxygens (O), and number of any other atoms (OTHER).

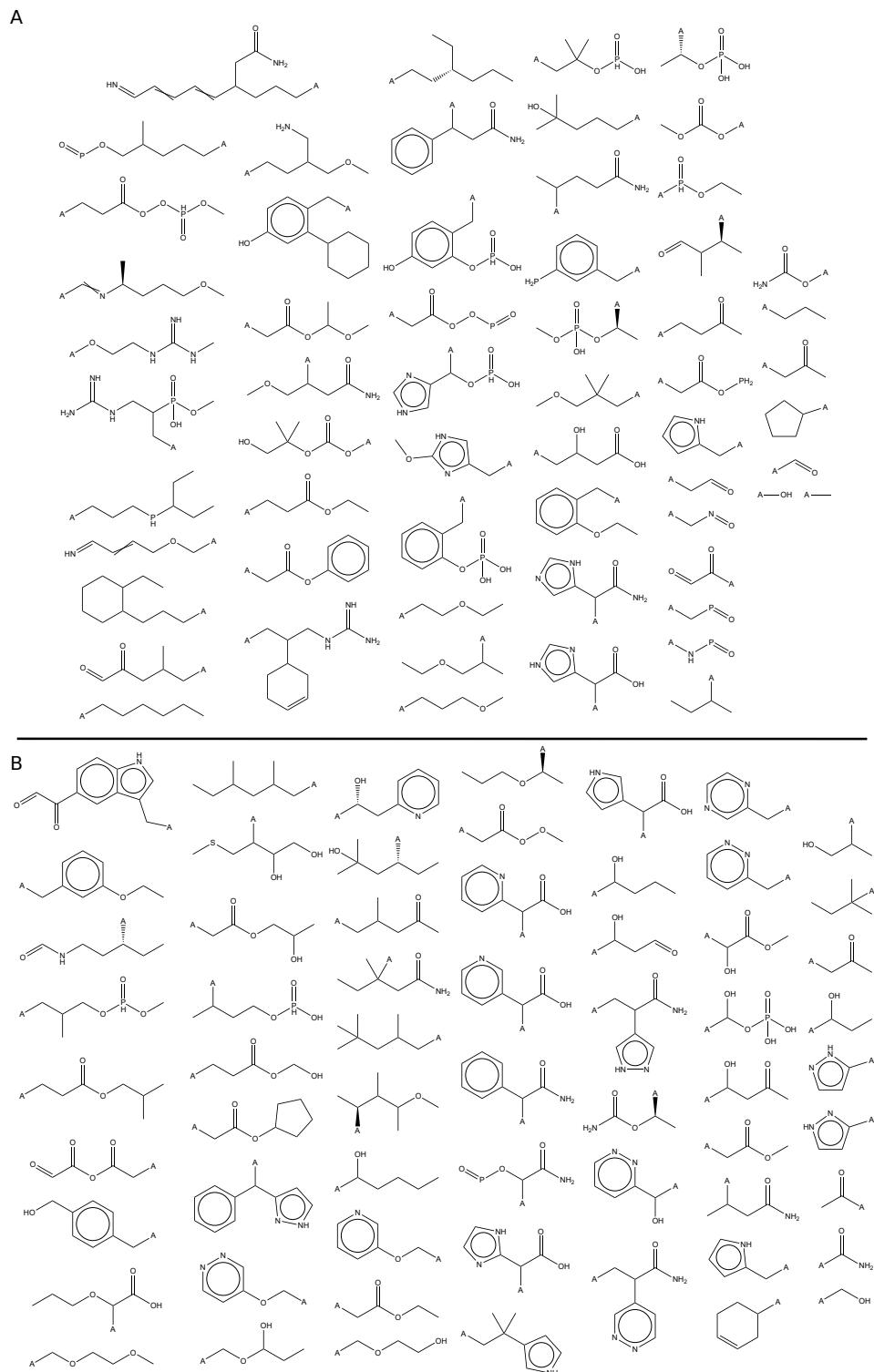


FIG. S3. Examples of side chains from a single training protein with unnatural amino acids and a single model generated protein with unnatural amino acids.

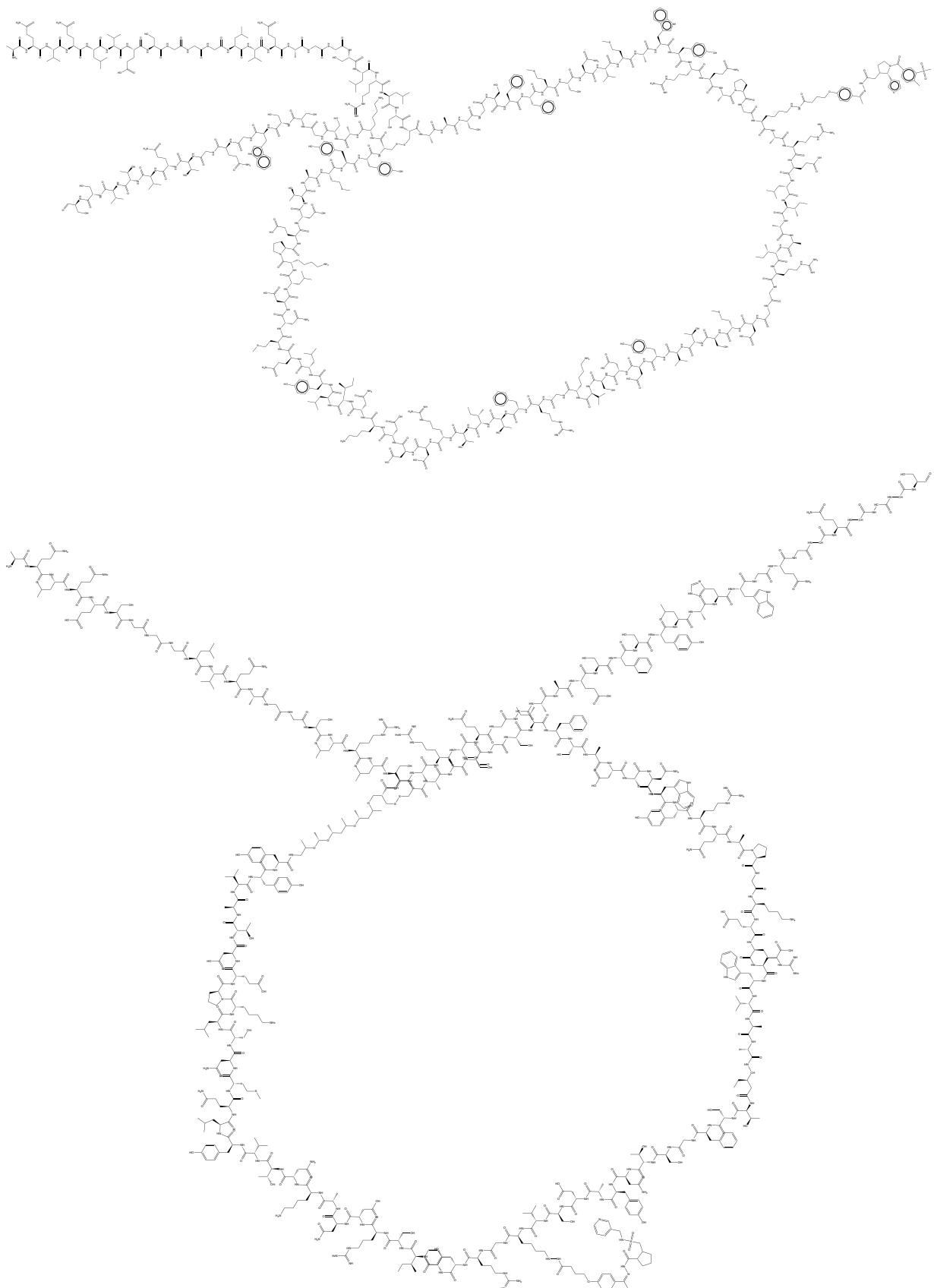


FIG. S4. **Antibody-Drug conjugates** Example antibody-drug conjugates from the training data (above) and one example produced by the language model (below) and plotted using ChemDraw as an atom-level graph.

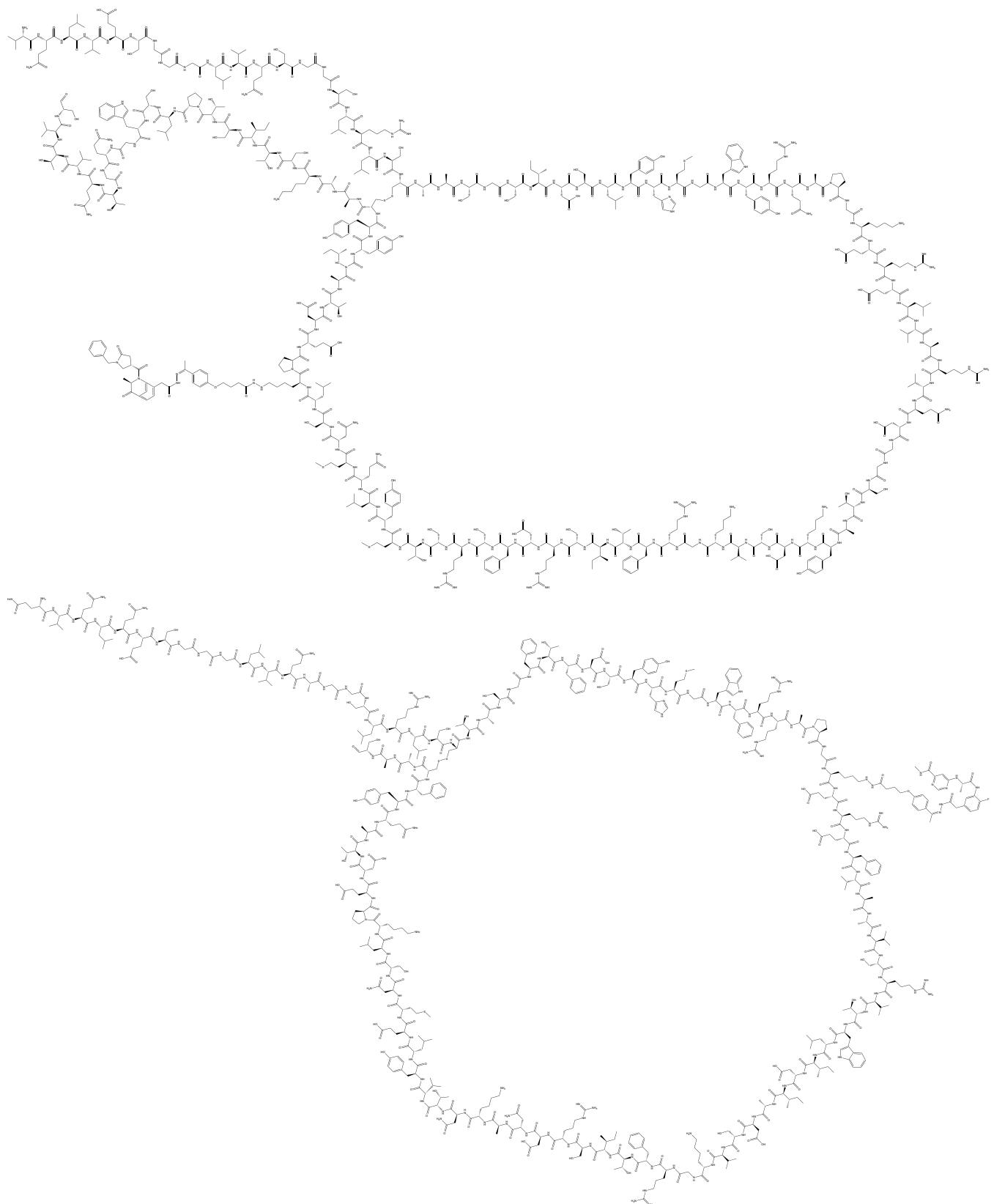


FIG. S5. Model Antibody-drug conjugates Example antibody-drug conjugates produced by the language model and plotted using ChemDraw as an atom-level graph.

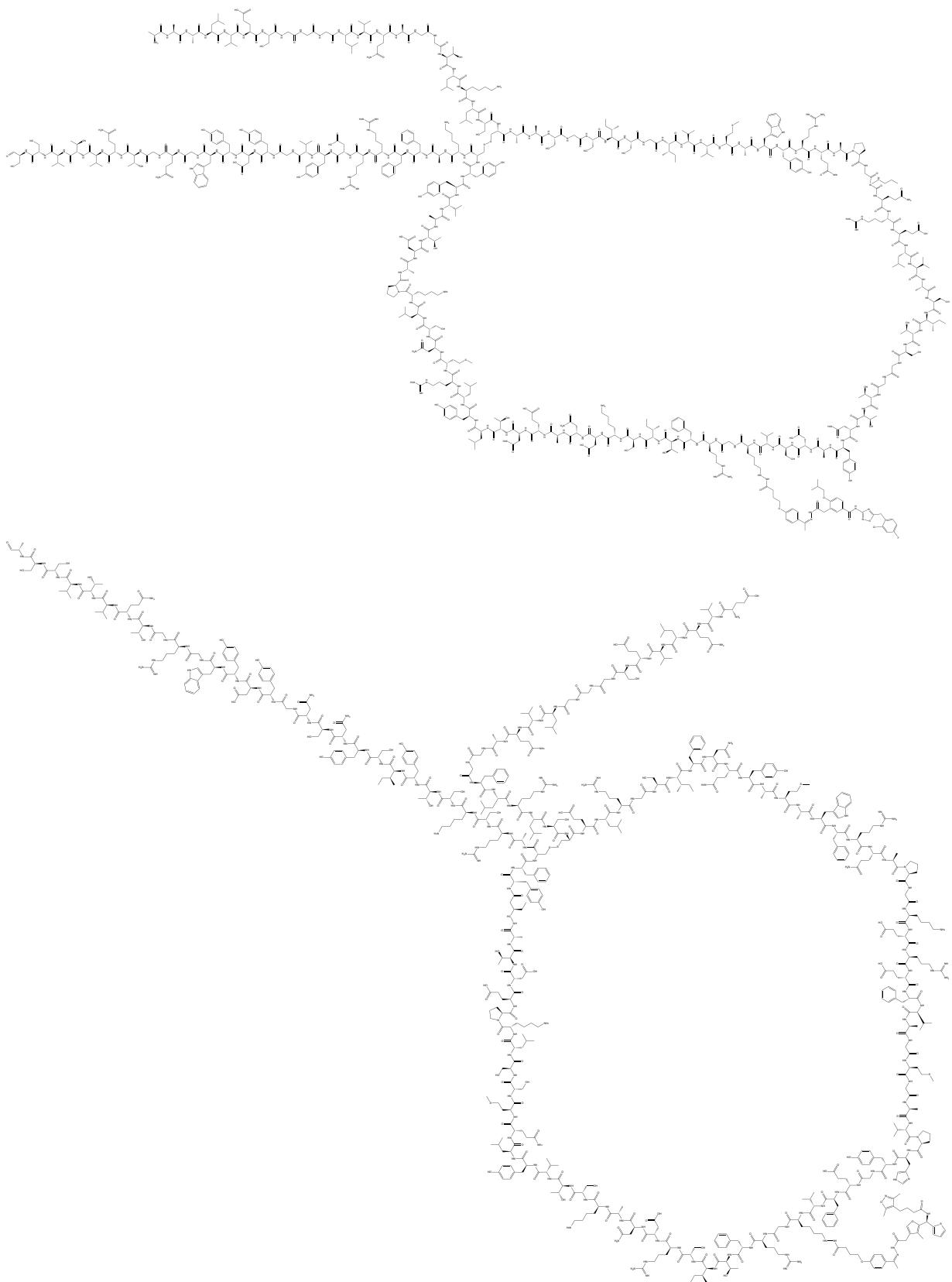


FIG. S6. **Model Antibody-drug conjugates** Example antibody-drug conjugates produced by the language model and plotted using ChemDraw as an atom-level graph.

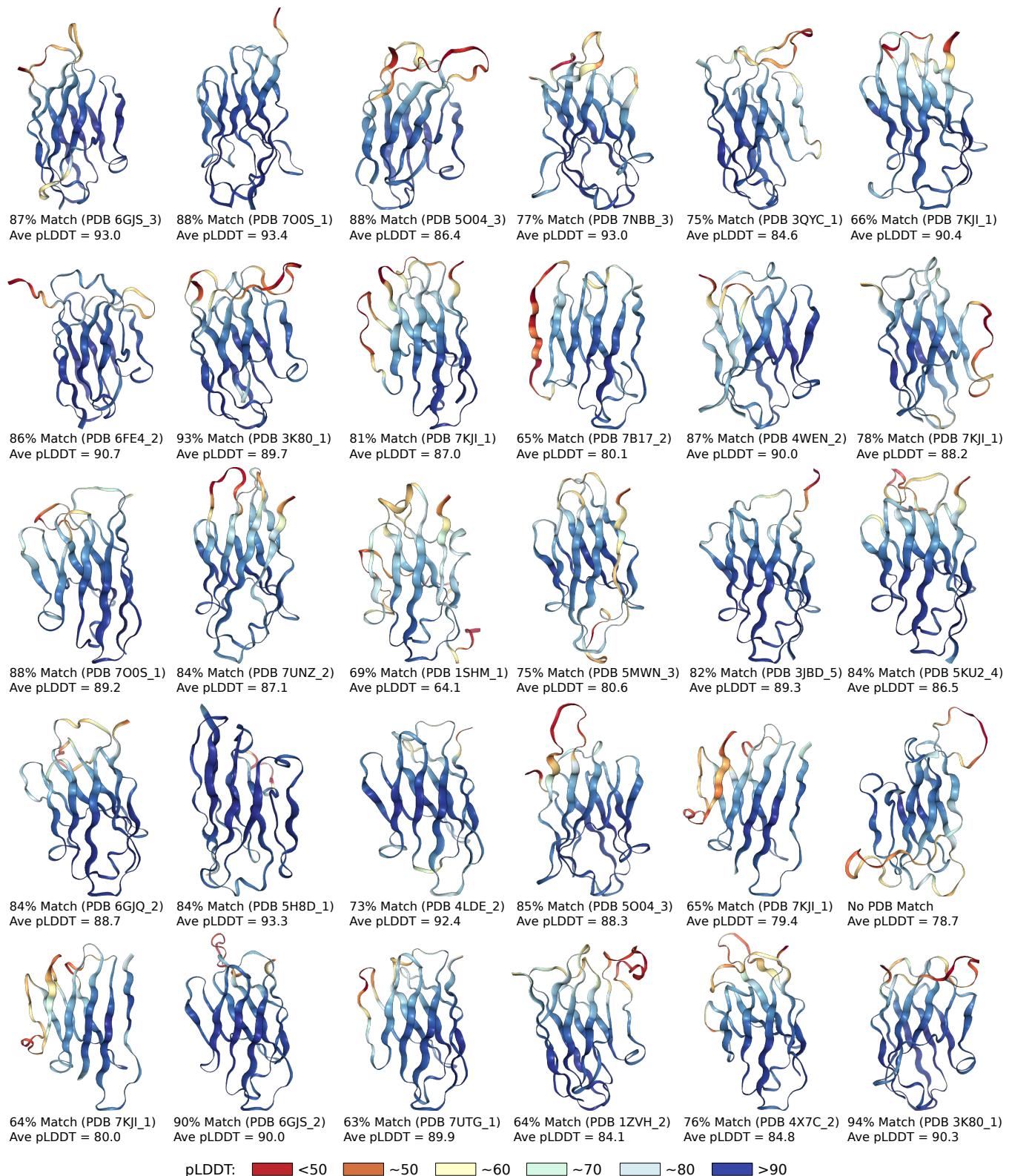


FIG. S7. Examples of single domain antibodies from samples of sdAb-drug conjugates (excluding warheads) generated by the model visualized by Alphafold and colored by pLDDT.

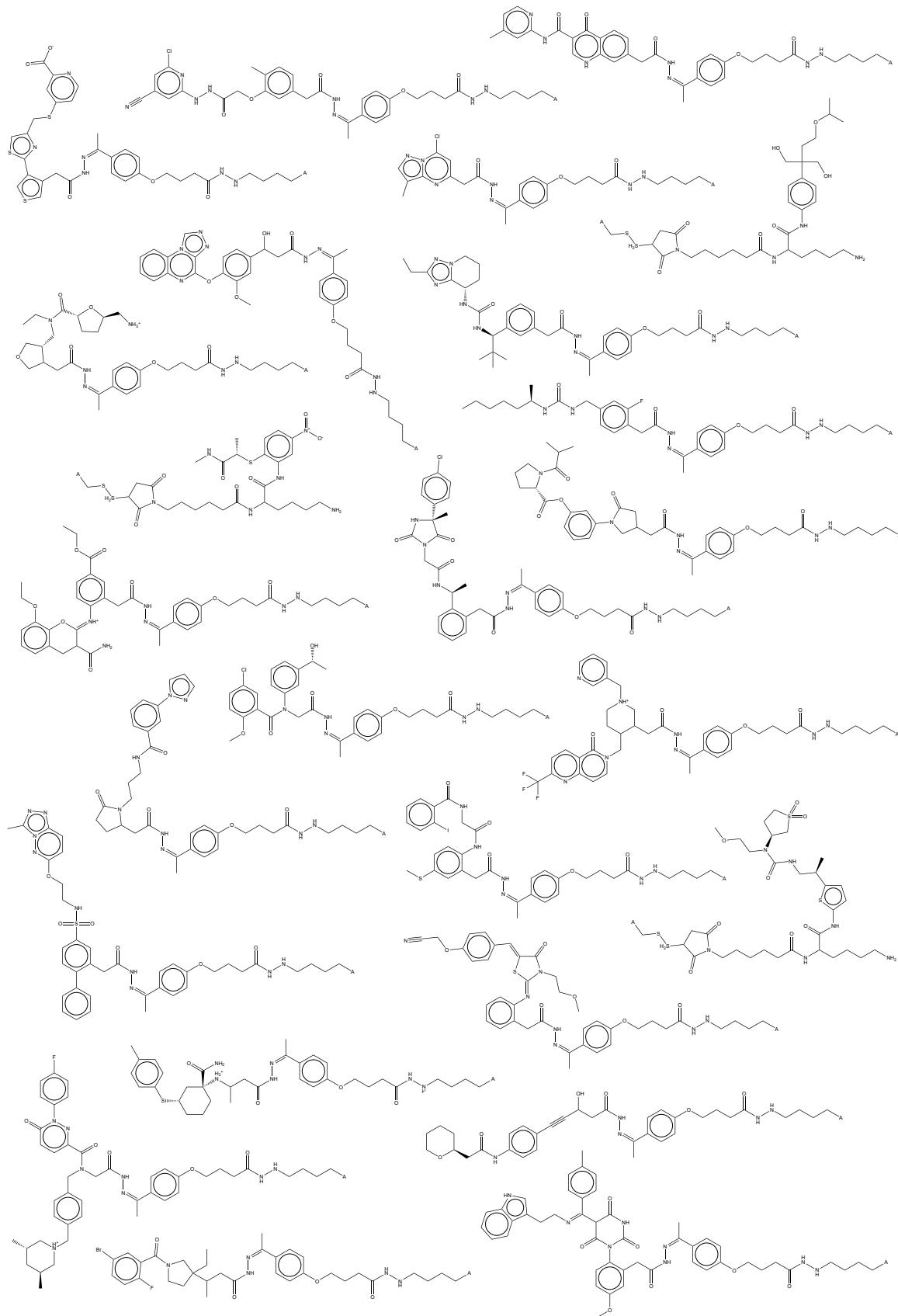


FIG. S8. Examples of detached warheads from training single domain antibody-drug conjugates.

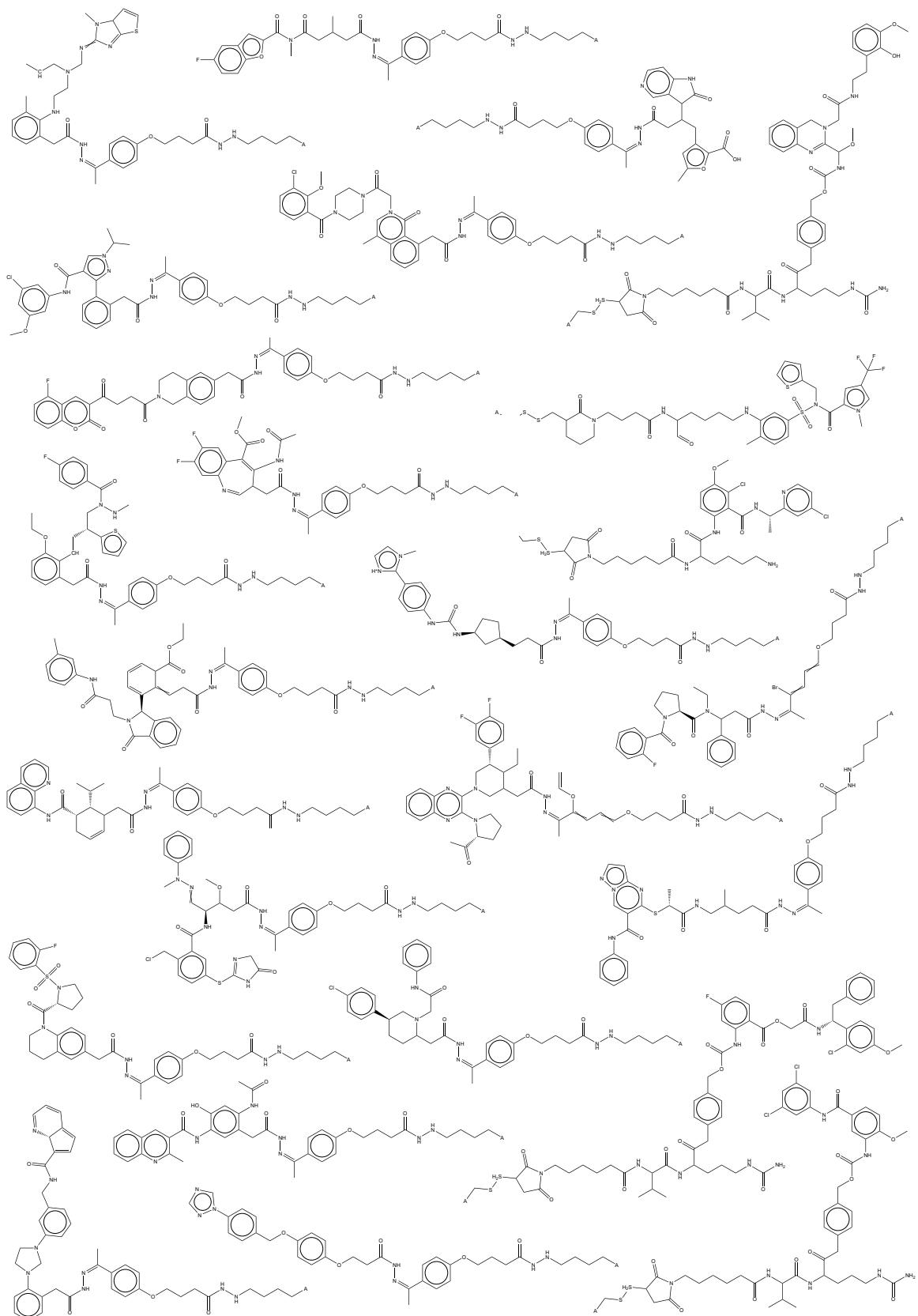


FIG. S9. Examples of detached warheads from model generated single domain antibody-drug conjugates.

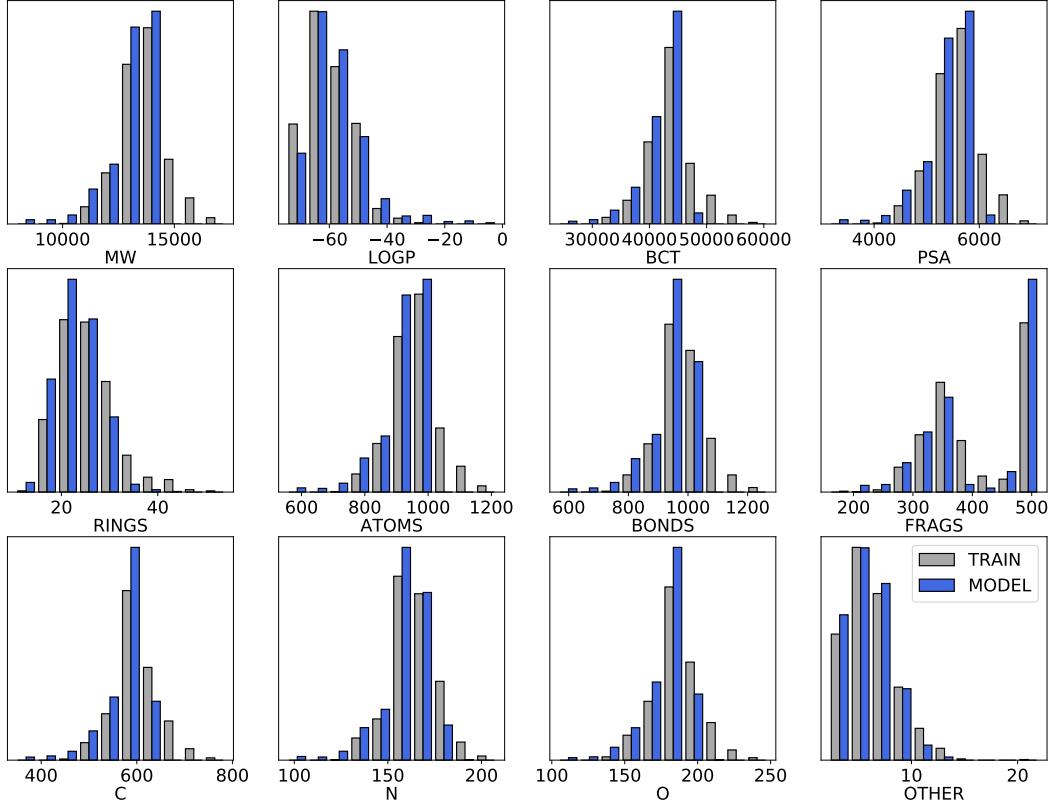


FIG. S10. Histograms of atom level properties for the antibody-drug-conjugates training data. These include exact molecular weight (MW), octanol-water partition coefficient (LogP) [17], molecular complexity (BCT), topological polar surface area (PSA), number of rings (RINGS), number of atoms (ATOMS), number of bonds (BONDS), number of fragments found by breaking up the molecule at rotatable bonds (FRAGS), number of carbons (C), number of nitrogens (N), number of oxygens (O), and number of any other atoms (OTHER)