
Integrating AI, automation and multiscale simulations for end-to-end design of phase-separating proteins

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Liquid-liquid phase separation (LLPS) is a fundamental cellular process that is
2 driven by self-assembly of intrinsically disordered proteins (IDPs), protein-RNA
3 complexes, or other bio-molecular systems which can form liquid droplets. Many
4 natural materials including silk, elastin, and gels are a result of LLPS and thus
5 rational design of such phase-separating peptides can have transformative impact,
6 from designing new biologically inspired materials (e.g., clothing) to self-
7 compartmentalized drug-delivery systems for biomedical applications. However,
8 given the intrinsic complexity in the rules governing LLPS, rational design of LLPS
9 undergoing peptides remains challenging. We posit that automation, foundation
10 models integrated with reinforcement learning approaches and multiscale molecular
11 simulations can drive the design of novel peptides that undergo LLPS. We describe
12 our progress towards the goal of end-to-end design of phase separating peptides
13 by summarizing current work at the Argonne National Laboratory’s Advanced
14 Photon Source 8ID-I beamline, where a robotic set up in the laboratory is enabled
15 via simulation and extensive testing of such bio-materials. Together, our approach
16 enables the design of novel bio-materials that can undergo phase separation under
17 diverse physiological conditions.

18 1 Introduction

19 Phase separation in biology is now being widely acknowledged as a fundamental mechanism of cellular control, including cellular compartmentalization as well as in various diseases such as cancer [Hyman et al. \[2014\]](#). More importantly, several naturally available proteins such as elastin [Rodríguez-Cabello et al. \[2018\]](#), silk [Lemetti et al. \[2022\]](#), [Parker et al. \[2019\]](#), and others are known to undergo phase separation which is likely to influence their overall stability and function inside of cells. Given that phase separation within such proteins is dependent on their (polymer) length, sequence specific linear (amino-acid) motifs, and other factors, a natural question is then in engineering novel constructs of such phase separating peptides/proteins that can possess specific properties [Hyman et al. \[2014\]](#).

27 Previous studies have examined how sequence composition and polymer length affect phase separation properties in elastin-like polypeptides (ELPs) [Christensen et al. \[2013\]](#). However, given the diversity of such sequences and the specific linear motifs that they need to phase separate under physiological conditions (e.g., ELPs utilize $-(VPXVG)_n-$ motif interspersed with other amino-acid sequences), the combinatorial complexity of the design space entails that an exhaustive evaluation of even a single class of phase separating peptides can be daunting, tedious, and error-prone. Furthermore, the discovery of new phase separating peptides/proteins with diverse mechanisms of self-assembly, there is a need to develop robust experimental and computational workflows that can probe and quantify how phase separation leads to different behaviors under diverse physiological conditions.

36 We posit that robotics and automa-
 37 tion within the laboratory integrated
 38 tightly with artificial intelligence (AI)
 39 methods, including generative models
 40 and reinforcement learning (RL)
 41 can provide an effective platform for
 42 not only characterizing phase separa-
 43 tion mechanisms, but also in design-
 44 ing novel peptides/proteins that un-
 45 dergo controlled phase separation un-
 46 der diverse conditions. As shown in
 47 Fig. 1, our automated design platform
 48 at Argonne National Laboratory inte-
 49 grates high-performance computing
 50 systems within the Argonne Leader-
 51 ship Computing Facility (ALCF) with
 52 the Advanced Photon Source (APS)
 53 beamline for characterizing phase sep-
 54 arating proteins using x-ray photon
 55 correlation spectroscopy (XPCS) and
 56 the Advanced Protein Characteriza-
 57 tion Facility (APCF) to clone, express,
 58 and purify protein samples at scale. At the heart of this *self driving lab* is a computational engine
 59 that consists of a suite of generative AI models that has been trained on diverse genome-scale data
 60 using large language models and fine-tuned on phase-separating protein databases. A RL approach
 61 is used to guide the precise modifications to the protein sequence that can predict specific phase
 62 separation properties. These are fed into a multiscale simulation framework that uses enhanced
 63 sampling techniques guided by AI approaches, namely, DeepDriveMD [Brace et al. [2022], Casalino
 64 et al. [2021]] to characterize molecular interactions that control phase separation. This approach lets
 65 us screen over 10^5 design candidates rapidly, while the APCF can automate the screening of 10^2 - 10^3
 66 sequence designs. The refined designs (about 10^2) are then characterized for phase separation at
 67 APS-8-ID-I beamline under physiological conditions and the observations are automatically ‘piped’
 68 through training the AI approaches (so that the design space can be constrained and conditioned
 69 appropriately). We provide an overview of progress in developing each of the areas highlighted.

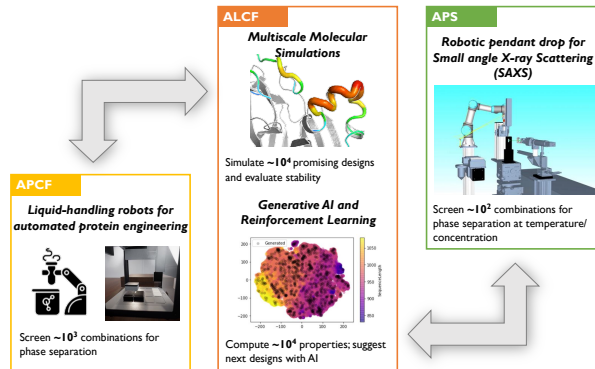


Figure 1: **An overview of our automated platform for designing phase-separating proteins.** Protein engineering is enabled via liquid handling robots enabling screening of 10^3 peptides. These are fed into simulation and AI workflows that automatically suggest new protein designs for subsequent rounds. Finally, a smaller set of protein designs (about 10^2) are characterized using X-ray scattering approaches.

70 2 AI-enabled phase-separating protein design and multiscale simulations

71 2.1 Reinforcement learning and generative sequence models

72 We formulate the design of a single sequence x using reinforcement learning where a policy is trained
 73 to optimize a specific objective [Sutton and Barto [2018], Silver et al. [2016]]. We initialize a policy
 74 $\pi = \rho$, where ρ is a pretrained generative language model providing the conditional probability
 75 distribution to predict the next tokens in the sequence. The initialized policy π is then fine-tuned
 76 using RL to perform the protein sequence-specific generation task. Combination of RL and language
 77 models have been successful in the past, where RL models were applied to fine-tune pre-trained
 78 language models for tasks such as text continuation with positive sentiment or physically descriptive
 79 language and summarization [Ziegler et al. [2019]]. In our work, ρ is obtained using GPT-NeoX [Black
 80 et al. [2022]] trained on diverse protein sequence datasets and fine-tuned on the phase-separating
 81 protein databases containing $\sim 6K$ sequences [You et al. [2020]].

82 We employ the proximal policy optimization (PPO) algorithm [Schulman et al. [2017]] as the RL
 83 model. The PPO policy guides the agent’s actions which in this case is to insert an amino-acid token
 84 from the sequence model vocabulary. The vocabulary consists of 21 amino-acid tokens and other
 85 special tokens as part of the tokenization process. In the initial set of experiments for generating
 86 novel sequences, the reward structure is simplified such that the reward structure benefits insertion of
 87 (valid) amino-acid tokens and penalized for adding special tokens, with a maximum length of 512
 88 tokens. Our experiments suggest (Fig. 2) that the RL training results in generating novel sequences.

89 2.2 Multiscale Simulations

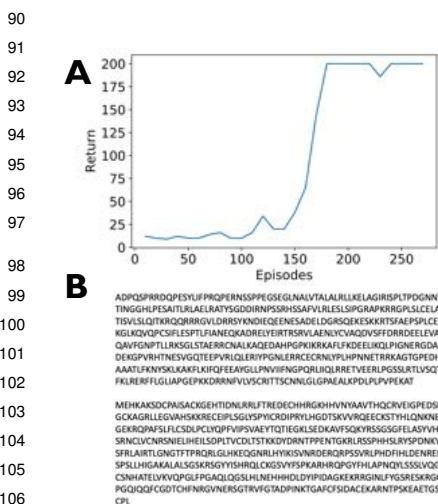


Figure 2: **RL-generated examples of phase-separating proteins.** (A) Returns (expected sum of rewards) are shown to improve with episodes suggesting that the PPO policy in combination with the GPT-NeoX model learns to generate protein sequences by guiding the agent’s actions to insert (valid) amino-acid tokens. (B) Examples of novel protein sequences generated after the training with a maximum sequence length of 512 tokens.

evaluated using a cutoff of 0.9nm. Particle-mesh Ewald summation was used to calculate the long-range electrostatic interactions with an Ewald error tolerance of 0.0005 and Hydrogen mass repartitioning to 1.5amu (to accelerate the integration). Preliminary benchmarking of the REMD give 930ns/day of simulation for systems having 10^5 atoms on 8 GPU (NVIDIA A100 cards).

We conducted large-scale replica exchange molecular dynamics (REMD) simulations of LLPS phase-separation in the generated peptides undergoing [Dignon et al. \[2019\]](#) to characterize the inter- and intra-molecular interactions that influence LLPS. We separately simulated (i) diffusion of the individual peptides through explicit solvent and (ii) closer-range interactions amongst multiple peptides (i.e., peptide aggregation) in explicit solvent.

The peptides consist of 70-150 amino acids and initially simulated in implicit solvent (see below) for 20ns to reach stable equilibrium conformations; equilibration of the RMSD and radius of gyration occurred within 20ns of simulation. Individual peptides were then simulated in an explicit solvent model (i.e., water and 150mM NaCl in a box providing ≥ 2 nm padding around the peptide). The peptide and explicit solvent contained roughly 10^5 atoms. Multiple peptides were similarly simulated in an explicit solvent model, pandas except they were arranged in a $3 \times 3 \times 3$ cuboid configuration with 10nm center-to-center distance between adjacent peptides. The multiple peptide systems contained roughly 10^6 atoms. These explicit solvent systems used the ff99sb.ILDN force field and TIP3P water model [Lindorff-Larsen et al. \[2010\]](#).

Replica exchange simulations were carried at 64 temperatures between 279.15 and 450 K. Each replica used the Langevin integrator with $1.0ps^{-1}$ collision rate and 0.004ps time-steps. The replicas were integrated for 2ps between each attempted exchange. The short-range electrostatic interactions and Lennard-Jones interactions were

123 3 Automated phase-separating protein engineering

The DNA fragments encoding the selected peptide repeats were generated by the overlap-extension rolling circle amplification (OERCA) method [Amiram et al. \[2011\]](#). The generated clones were sequenced and those with repeats of 20 or more peptides were selected for characterization. The N-terminal His₆-tagged proteins were purified using immobilized metal-affinity chromatography and used without the removal of the purification tag. The phase transition of the proteins was measured by either monitoring absorbance at 350nm in a plate reader or monitoring interaction with a fluorescent dye via a real-time PCR detection system (Fig. 3). The phase transition temperature was lower for polypeptides with fewer repeats. Similarly, lower protein concentration, the addition of NaCl or PEG-8000 to the solution also resulted in lower phase transition temperatures. A fully automated system, driven by a Python API, can be used to assemble the various combinations or polypeptides and additives. Current efforts focus on the development of a fully automated closed-loop system capable not only measuring phase transition of a given input

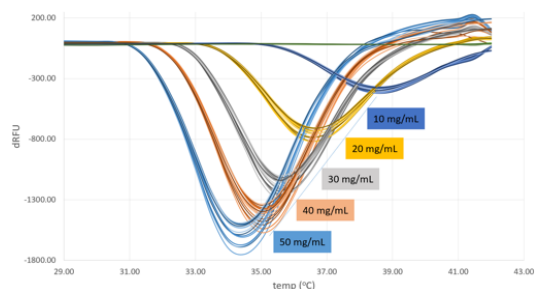
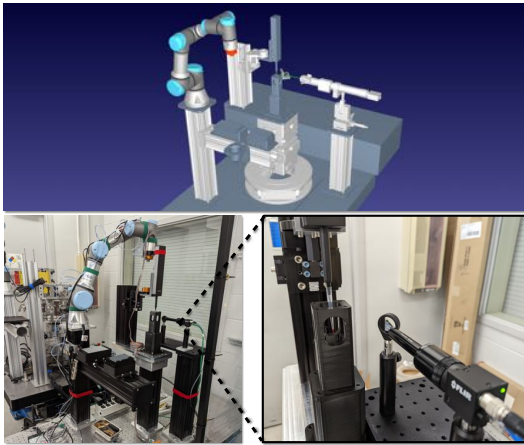


Figure 3: Example phase transition measurement at the Advanced Protein Characterization Facility. The phase transition temperature was lower for polypeptides with fewer repeats. Similarly, lower protein concentration, the addition of NaCl or PEG-8000 to the solution also resulted in lower phase transition temperatures. A fully automated system, driven by a Python API, can be used to assemble the various combinations or polypeptides and additives. Current efforts focus on the development of a fully automated closed-loop system capable not only measuring phase transition of a given input

143 sample without human intervention, but eventually also carry out the synthesis, cloning, and protein
144 purification steps of the workflow.

145 4 Robotic pendant drop enabled small-angle scattering experiments

146 Robotic pendant drop setup was developed in the adjacent chemistry laboratory of beamline 8-ID-I of
147 Advanced Photon Source and robot programs were implemented on a simulation software [RoboDK](#)
148 [\[2022\]](#). To perform the pendant drop experiments [Bera and Antonio \[2016\]](#), UR3e collaborative robot
149 arm from Universal Robots was utilized as the liquid handling robot. UR3e robot was ideal to operate
150 in tight workspaces such as the beamline, due to its compactness and small footprint. In order to
151 create precise droplets and eliminate the vibration factor, the experimental setup was designed with a
152 pipette docking location. Furthermore, a tool changer (ATI QC-11) was attached to both the end joint
153 of the robot and an Opentrons Single Channel P300 GEN2 electronic pipette. The tool changer lets
154 us lock and unlock its Master and Tool sides with air compression to pick and place the pipette. The
155 pipette was driven to accurately control the volume of the liquid aspirated and dispensed inside the
156 tips and ejecting the tips when needed. An optical microscope was placed by the sample location to
157 provide live video feed of the sample via the reflection of the 45-degree mirror. The mirror has a 1
158 mm through-hole at its center, allowing x-ray beam to pass through so that optical inspection and
159 x-ray measurements can be performed simultaneously.



160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177 Figure 4: Robotic Setup in The Chemistry Laboratory of Beamline 8-ID-I of APS/Argonne.

178
179
180 respectively. Fig. 4 shows RoboDK is executing the program on the UR3e, and the robot is performing
181 liquid handing and sample exchange.

182 5 Summary

183 We have highlighted our progress in developing a self-driving laboratory for designing phase separ-
184 ating proteins. Our approach uses robotics integrated in a functioning beamline (to characterize
185 size and dynamics of phase separation) with AI/ML techniques and high-throughput molecular
186 simulations. The approach also highlighted some important lessons that we learned, including the
187 challenges involved in integrating diverse robotic systems and how such ecosystems of commercial
188 off-the-shelf robotic systems integration can be carried out across at a user facility. Further, it also
189 highlighted the importance of building robust, automated workflow systems that can be used to enable
190 high-throughput bio-materials characterization. Finally, it also helped us reduce the time-to-solution
191 for design cycles of phase separating proteins from several months to about weeks – thus allowing a
192 much rapid exploration of the design space of such materials.

193 References

194 Miriam Amiram, Felipe Garcia Quiroz, Daniel J. Callahan, and Ashutosh Chilkoti. A highly parallel
195 method for synthesizing dna repeats enables the discovery of ‘smart’ protein polymers. *Nature*

- 196 *Materials*, 10(2):141–148, 2011. doi: 10.1038/nmat2942. URL <https://doi.org/10.1038/nmat2942>.
197
- 198 Mrinal K Bera and Mark R Antonio. Crystallization of keggins heteropolyanions via a Two-Step
199 process in aqueous solutions. *J. Am. Chem. Soc.*, 138(23):7282–7288, June 2016.
- 200 Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He,
201 Connor Leahy, Kyle McDonnell, Jason Phang, et al. Gpt-neox-20b: An open-source autoregressive
202 language model. *arXiv preprint arXiv:2204.06745*, 2022.
- 203 Alexander Brace, Igor Yakushin, Heng Ma, Anda Trifan, Todd Munson, Ian Foster, Arvind Ra-
204 manathan, Hyungro Lee, Matteo Turilli, and Shantenu Jha. Coupling streaming ai and hpc
205 ensembles to achieve 100–1000× faster biomolecular simulations. In *2022 IEEE Interna-
206 tional Parallel and Distributed Processing Symposium (IPDPS)*, pages 806–816, 2022. doi:
207 10.1109/IPDPS53621.2022.00083.
- 208 Lorenzo Casalino, Abigail C Dommer, Zied Gaieb, Emilia P Barros, Terra Sztain, Surl-Hee Ahn,
209 Anda Trifan, Alexander Brace, Anthony T Bogetti, Austin Clyde, Heng Ma, Hyungro Lee, Matteo
210 Turilli, Syma Khalid, Lillian T Chong, Carlos Simmerling, David J Hardy, Julio DC Maia,
211 James C Phillips, Thorsten Kurth, Abraham C Stern, Lei Huang, John D McCalpin, Mahidhar
212 Tatineni, Tom Gibbs, John E Stone, Shantenu Jha, Arvind Ramanathan, and Rommie E Amaro.
213 Ai-driven multiscale simulations illuminate mechanisms of sars-cov-2 spike dynamics. *The
214 International Journal of High Performance Computing Applications*, 35(5):432–451, 2021. doi:
215 10.1177/10943420211006452. URL <https://doi.org/10.1177/10943420211006452>.
- 216 Trine Christensen, Wafa Hassouneh, Kimberley Trabbic-Carlson, and Ashutosh Chilkoti. Predicting
217 transition temperatures of elastin-like polypeptide fusion proteins. *Biomacromolecules*, 14(5):
218 1514–1519, 05 2013. doi: 10.1021/bm400167h. URL <https://doi.org/10.1021/bm400167h>.
- 219 Gregory L Dignon, Wenwei Zheng, and Jeetain Mittal. Simulation methods for liquid–liquid
220 phase separation of disordered proteins. *Current Opinion in Chemical Engineering*, 23:92–98,
221 2019. ISSN 2211-3398. doi: <https://doi.org/10.1016/j.coche.2019.03.004>. URL <https://www.sciencedirect.com/science/article/pii/S2211339818300807>. *Frontiers of Chemical
222 Engineering: Molecular Modeling*.
223
- 224 Anthony A. Hyman, Christoph A. Weber, and Frank Jülicher. Liquid-liquid phase sep-
225 aration in biology. *Annual Review of Cell and Developmental Biology*, 30(1):39–58,
226 2014. doi: 10.1146/annurev-cellbio-100913-013325. URL <https://doi.org/10.1146/annurev-cellbio-100913-013325>. PMID: 25288112.
227
- 228 Martin R Kraimer, Janet B Anderson, Andrew N Johnson, W Eric Norum, Jeffrey O Hill, Ralph
229 Lange, Benjamin Franksen, and Peter Denison. Epics application developer’s guide. *EPICS Base
230 Release*, 3(11):1–243, 2012.
- 231 Laura Lemetti, Alberto Scacchi, Yin Yin, Mengjie Shen, Markus B. Linder, Maria Sammalkorpi, and
232 A. Sesilja Aranko. Liquid–liquid phase separation and assembly of silk-like proteins is dependent
233 on the polymer length. *Biomacromolecules*, 23(8):3142–3153, 08 2022. doi: 10.1021/acs.biomac.
234 2c00179. URL <https://doi.org/10.1021/acs.biomac.2c00179>.
- 235 Kresten Lindorff-Larsen, Stefano Piana, Kim Palmo, Paul Maragakis, John L. Klepeis, Ron O.
236 Dror, and David E. Shaw. Improved side-chain torsion potentials for the amber ff99sb protein
237 force field. *Proteins: Structure, Function, and Bioinformatics*, 78(8):1950–1958, 2010. doi:
238 <https://doi.org/10.1002/prot.22711>. URL [https://onlinelibrary.wiley.com/doi/abs/10.
239 1002/prot.22711](https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.22711).
- 240 Rachael N. Parker, Wenyao A. Wu, Tina B. McKay, Qiaobing Xu, and David L. Kaplan. Design of
241 silk-elastin-like protein nanoparticle systems with mucoadhesive properties. *Journal of Functional
242 Biomaterials*, 10(4), 2019. ISSN 2079-4983. doi: 10.3390/jfb10040049. URL [https://www.
243 mdpi.com/2079-4983/10/4/49](https://www.mdpi.com/2079-4983/10/4/49).
- 244 RoboDK. Robodk: Simulate robot applications., 2022. URL <https://robodk.com/>.

- 245 Jose C. Rodríguez-Cabello, Israel González de Torre, Sergio Acosta, Soraya Salinas, and Marcos
246 Herrero. 4 - elastin-like proteins: Molecular design for self-assembling. In Helena S. Azevedo
247 and Ricardo M.P. da Silva, editors, *Self-assembling Biomaterials*, Woodhead Publishing Series
248 in Biomaterials, pages 49–78. Woodhead Publishing, 2018. ISBN 978-0-08-102015-9. doi:
249 <https://doi.org/10.1016/B978-0-08-102015-9.00004-6>. URL [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/B9780081020159000046)
250 [com/science/article/pii/B9780081020159000046](https://www.sciencedirect.com/science/article/pii/B9780081020159000046).
- 251 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
252 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 253 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, and Laurent Sifre. George van den
254 driessche et al., mastering the game of go with deep neural networks and tree search. *Nature*, 529:
255 7587, 2016.
- 256 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 257 Kaiqiang You, Qi Huang, Chunyu Yu, Boyan Shen, Cristoffer Sevilla, Minglei Shi, Henning Herm-
258 jakob, Yang Chen, and Tingting Li. Phasepdb: a database of liquid–liquid phase separation related
259 proteins. *Nucleic acids research*, 48(D1):D354–D359, 2020.
- 260 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
261 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*
262 *preprint arXiv:1909.08593*, 2019.