

REINFORCEMENT LEARNING FOR BLACK-BOX OBJECTIVES IN LOGIT-BASED PROTEIN HALLUCINATION

Anonymous authors
Paper under double-blind review

ABSTRACT

Gradient-based “hallucination” methods such as Germinal and BindCraft enable efficient protein binder design by optimizing a continuous relaxation of the sequence (a logit matrix) using gradients from differentiable structure predictors and protein language models. However, many practically useful objectives for ranking or filtering designs are non-differentiable with respect to sequence (e.g., external confidence metrics from AlphaFold3 or Chai, experimental readouts, or arbitrary black-box scores), preventing direct backpropagation through the objective. We extend the Germinal pipeline with an *optional* policy-gradient update on the sequence logits, enabling direct optimization of black-box rewards while preserving Germinal’s differentiable optimization backbone. Our implementation reuses Germinal’s existing filter metrics as modular reward components and supports both Chai-1 and AlphaFold3 backends for reward evaluation. On a nanobody design task against PD-L1, adding a small policy-gradient term during the *high-softness* portion of Germinal’s optimization yields a $2.3\times$ improvement in acceptance rate among processed trajectories, from 0.161 ± 0.020 to 0.370 ± 0.055 (mean \pm std over six seeds), while maintaining comparable confidence metrics for accepted designs. This suggests that combining policy-gradient-based black-box optimization helps improve design success rates, potentially improving downstream wet-lab metrics.

1 INTRODUCTION

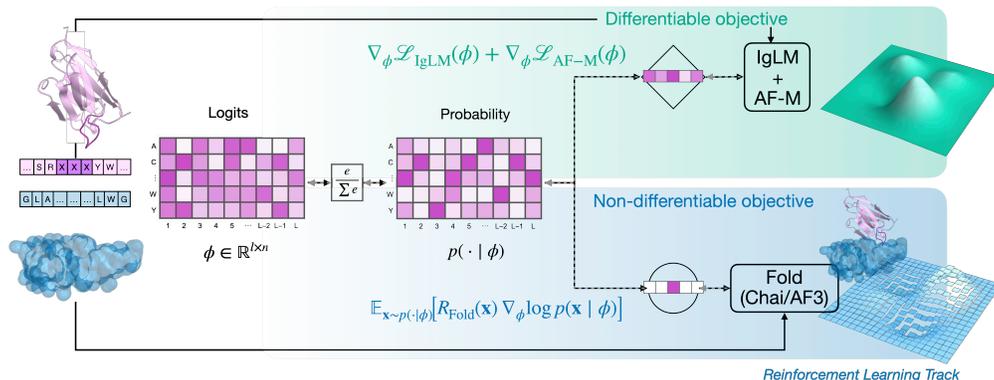


Figure 1: Overview of the RL-augmented logit-optimization pipeline. Logits are updated by combining differentiable hallucination gradients with a policy-gradient term computed from black-box rewards.

Designing functional protein binders requires searching a large discrete sequence space under structural and interface constraints. Recent progress has come from gradient-based design pipelines that optimize sequence logits through differentiable structure models and sequence priors Frank et al. (2024); Pacesa et al. (2025); Mille-Fragoso et al. (2025). These approaches are effective, but they are naturally limited to objectives that can be backpropagated.

Two families dominate current practice. Hallucination-based methods optimize sequences directly against structure-model losses and can produce strong sequence–structure consistency Frank et al. (2024); Pacesa et al. (2025); Mille-Fragoso et al. (2025). Diffusion-based pipelines sample candidate

054 structures quickly and then design sequences conditioned on those structures [Bennett et al. \(2024\)](#);
 055 [Watson et al. \(2023\)](#); [Gruver et al. \(2023\)](#). Both families still rely heavily on downstream filtering and
 056 cannot directly optimize black-box success criteria. For a larger discussion of related works, please
 057 see [A.1](#).

058 For example, confidence metrics from external structure predictors like AlphaFold3 [Abramson
 059 et al. \(2024\)](#) are often useful for ranking candidates, but are not differentiable with respect to
 060 sequence. Reinforcement learning offers a natural route to optimize such objectives without requiring
 061 differentiability [da Silva et al. \(2023\)](#); [Ding et al. \(2024\)](#); [Lee et al. \(2025\)](#).

062 We propose a minimal, modular RL augmentation for gradient-based protein hallucination pipelines.
 063 Our method adds an optional policy-gradient term to the existing logit optimization loop, supports
 064 rewards from either internal design metrics or external predictors (e.g., Chai-1 or AlphaFold3 [Chai
 065 Discovery \(2024\)](#); [Abramson et al. \(2024\)](#)), and can be activated only in selected optimization regions
 066 (e.g., late logits / softmax) to control variance and compute. When RL is disabled, the pipeline is
 067 reduced to the baseline hallucination method.

068 The central question is practical: *can sparse RL updates improve end-to-end design yield without
 069 destabilizing the core hallucination workflow?* We evaluate this question on PD-L1 nanobody design
 070 and report acceptance-rate behavior across multiple independent seeds.

072 2 METHOD

073 RALO is implemented as an *optional extension* to Germinal [Mille-Fragoso et al. \(2025\)](#). We retain
 074 Germinal’s core optimization loop—gradient-based hallucination on a logit parameterization of
 075 sequences—and add a policy-gradient term that can optimize *black-box* objectives that do not admit
 076 backpropagation. Figure 1 provides a high-level overview.

078 2.1 BACKGROUND: GERMINAL OPTIMIZES A LOGIT DISTRIBUTION

079 Let $\mathbf{x} = (x_1, \dots, x_N)$ be an amino-acid sequence of length N over an alphabet of size M (typically
 080 $M = 20$). Germinal maintains logits $\phi \in \mathbb{R}^{N \times M}$ and defines an independent categorical distribution
 081 using Eq. (1). Gradients from a structure prediction loss (AlphaFold2-Multimer [Evans et al. \(2022\)](#))
 082 and sequence prior (IgLM [Shuai et al. \(2023\)](#)) are merged (e.g., with PCGrad [Yu et al. \(2020\)](#)) and
 083 applied to the logits during a three-phase schedule: *Phase 1 (logits)*, *Phase 2 (softmax)*, and *Phase 3*
 084 (*semi-greedy*).

086 2.2 POLICY-GRADIENT TERM FOR BLACK-BOX REWARDS

087 Many useful scoring functions for binder design are not differentiable with respect to ϕ , including
 088 confidence scores produced by external structure predictors (Chai-1, AlphaFold3 [Chai Discovery
 089 \(2024\)](#); [Abramson et al. \(2024\)](#)), or experimental measurements. To enable optimization under such
 090 objectives, we add an *optional* policy-gradient update. At selected iterations, we sample K sequences
 091 $\mathbf{x}^{(k)} \sim p(\cdot | \phi)$, evaluate a black-box reward $R(\mathbf{x}^{(k)})$, and apply the REINFORCE [Williams \(1992\)](#)
 092 estimator:

$$093 \nabla_{\phi} \mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [R(\mathbf{x})] \approx \frac{1}{K} \sum_{k=1}^K (R(\mathbf{x}^{(k)}) - b) \nabla_{\phi} \log p(\mathbf{x}^{(k)} | \phi),$$

095 where b is a baseline (we use a per-update batch mean) to reduce variance. In practice, we implement
 096 this by minimizing the surrogate loss

$$097 \mathcal{L}_{\text{RL}}(\phi) = -\frac{1}{K} \sum_{k=1}^K (\tilde{R}^{(k)} - b) \log p(\mathbf{x}^{(k)} | \phi),$$

098 where $b = \frac{1}{K} \sum_k \tilde{R}^{(k)}$ is the per-update baseline. For the setting used in this paper, the scalar reward
 099 is clipped as

$$100 \tilde{R}^{(k)} = \text{clip}(R(\mathbf{x}^{(k)}); -c, c),$$

101 where the clipping threshold c is a tunable hyperparameter. When reward variance within an update
 102 is effectively zero ($\text{std}(R) \leq \epsilon$, with ϵ a small constant), the implementation skips the RL update to
 103 avoid adding numerical noise. The total logit update combines the differentiable Germinal gradients
 104 with the RL term:

$$105 g_{\text{total}} = g_{\text{AF}} + g_{\text{IgLM}} + \lambda_{\text{RL}} g_{\text{RL}},$$

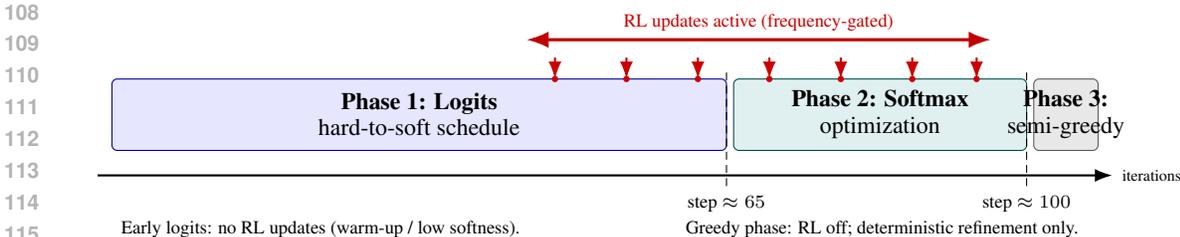


Figure 2: Temporal view of Germinal phases and where the optional RL gradient is injected. In the RALO configuration, RL is applied only in late logits and softmax steps, at frequency-gated update events. The RL gradient is added to the logit-gradient tensor and merged with standard Germinal gradients before each optimizer step.

with a tunable coefficient λ_{RL} . When $\lambda_{\text{RL}} = 0$ (or the feature flag is disabled), behavior is identical to the unmodified Germinal pipeline.

2.3 REWARD CONSTRUCTION FROM GERMINAL METRICS

We define $R(\mathbf{x})$ as a weighted sum of modular reward components: $R(\mathbf{x}) = \sum_i w_i r_i(\mathbf{x})$. To maximize reuse and avoid backend-specific code paths, each r_i wraps an existing Germinal Mille-Fragoso et al. (2025) metric. In particular, we support: (i) *internal* metrics from Germinal (or optional per-sample AlphaFold Evans et al. (2022) reevaluation), and (ii) *external* metrics computed by an external predictor (Chai-1 or AlphaFold3 Chai Discovery (2024); Abramson et al. (2024)), using the same backend selected by the run configuration.

Margin Mode. To increase filter pass rates, we introduce a margin reward. Given a filter threshold τ and a metric $m(\mathbf{x})$, we define a signed margin:

$$r^{\text{margin}}(\mathbf{x}) = \begin{cases} m(\mathbf{x}) - \tau, & \text{for filters of the form } m(\mathbf{x}) \geq \tau, \\ \tau - m(\mathbf{x}), & \text{for filters of the form } m(\mathbf{x}) \leq \tau. \end{cases}$$

Positive margin indicates the constraint is satisfied; negative margin indicates violation. This simple transformation makes the RL term “push” designs toward passing the same constraints used downstream.

2.4 SCHEDULING: APPLY RL WHEN THE DISTRIBUTION HAS VARIANCE

Policy-gradient updates require reward variance under $p(\cdot | \phi)$. If the logit distribution collapses to (nearly) a point mass, or if sampling is deterministic, then the sampled rewards become (nearly) constant and the REINFORCE update becomes ineffective. We therefore introduce simple scheduling controls: (i) a warm-up step, (ii) a minimum “softness” threshold, (iii) an update frequency, and (iv) a sampling strategy. Figure 2 provides a temporal view of these controls across Germinal’s three phases and highlights where RL updates are active. Here “softness” refers to Germinal’s scalar `soft` knob that interpolates between a hard, nearly deterministic sequence representation and a fully soft categorical distribution over amino acids. In the experiments in Section 3, we restrict RL updates to the late steps of the logits phase and the softmax phase; this controls compute and helps avoid no-variance updates.

3 EXPERIMENTS

We evaluate RALO as an *extension* of the Germinal pipeline Mille-Fragoso et al. (2025) on a nanobody binder design task against PD-L1. All runs use the same Germinal backbone, which attempts 50 trajectories and uses Germinal’s standard three-phase schedule (logits \rightarrow softmax \rightarrow semi-greedy). We add a small RL term in the high-softness region of optimization and test its effect on final filter acceptance rates (see Appendix A.4 for filter details). Our baseline is the unmodified Germinal pipeline, which we compare to an RL-augmented configuration that uses sparse policy-gradient updates with rewards computed by Chai-1 Chai Discovery (2024).

RL-augmented. We enable the RL term in the *high-softness* region of optimization (when logits are in their raw and softmax phases) and gate RL contributions until after a warm-up period. We compute

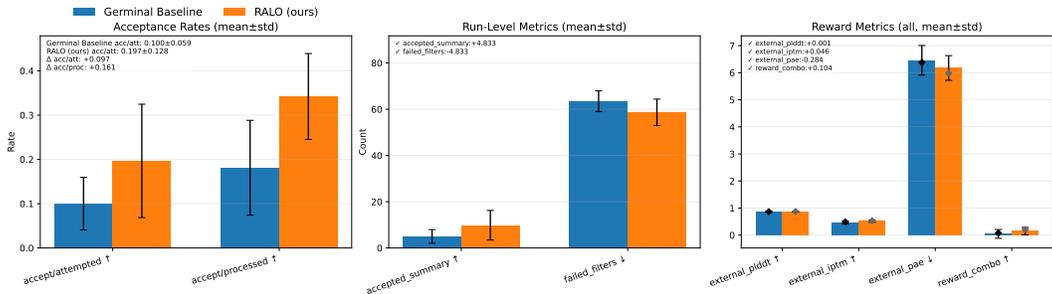


Figure 3: Baseline Germinal vs. RL-augmented Germinal (RALO) on PD-L1 nanobody design. Left: acceptance rates (mean±std across six seeds). Middle: counts of accepted and failed designs over entire runs (mean±std across six seeds). Right: accepted-set metrics (mean±std across seeds of per-run accepted-set means).

a margin reward as described above over the external confidence metrics already used by Germinal filters: pLDDT, ipTM, and PAE, all computed by Chai-1 with weights $w = \{+1, +1, -0.2\}$ and reward clipping $c = 2.0$. We enable sparse, small RL updates ($\lambda_{RL} = 0.02$) every 20 optimization steps. See Appendix A.3 for details.

Protocol and metrics. For each condition (baseline and RL), we run 6 independent seeds (6 separate executions). We report mean±std over seeds. We track two acceptance rates: (i) *accept/attempted*, the fraction of attempted trajectories that yield an accepted design, and (ii) *accept/processed*, the fraction of trajectories that reach the final evaluation stage that are accepted.

Results. Across the seeds, adding a small RL term in the high-softness region increases acceptance rates. The baseline achieves *accept/attempted* 0.100 ± 0.059 and *accept/processed* 0.181 ± 0.107 . The RL-augmented configuration achieves *accept/attempted* 0.197 ± 0.128 and *accept/processed* 0.342 ± 0.097 . The higher variance in *accept/attempted* results stems from the stochastic sampling performed by policy gradients; this variance tightens in *accept/processed* as RL contributions cease in the final optimization stage. These results suggest policy gradients enhance sequence space exploration, thus improving final performance. Figure 3 summarizes both acceptance rates, design counts, and reward metrics (Chai-1 confidence scores) to verify that improvements in acceptance rates are accompanied by improvements in the underlying reward metrics. The RL-augmented configuration achieves higher acceptance rates by improving the underlying reward metrics, which are directly related to the final filters.

4 CONCLUSION

We propose RALO, an extension to the Germinal protein design pipeline to enable optimization of important downstream signals ordinarily reserved for filtering of designs. We leverage policy gradients to optimize these objectives during the raw logits and softmax phases.

Our results show that integrating policy gradient methods for black-box optimization accelerates antibody design: across run seeds, RALO significantly increases the number of passing designs.

A limitation of our results is that our evaluation is limited to a single target (PD-L1) and nanobody scaffold. Future work should expand evaluation to diverse targets and scaffolds.

Future work will also explore greater algorithmic changes in more depth, study the effect of adding a policy gradient optimization step on generated designs, expand the *in silico* benchmark suite, and include wet lab validation of selected designs.

REFERENCES

Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein, David A. Evans, Chia-Chun Hung, Michael O’Neill, David Reiman, Kathryn Tunyasuvunakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles Congreve, Alexander I. Cowen-Rivers, Andrew

- 216 Cowie, Michael Figurnov, Fabian B. Fuchs, Hannah Gladman, Rishub Jain, Yousuf A. Khan,
217 Caroline M. R. Low, Kuba Perlin, Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian
218 Stecula, Ashok Thillaisundaram, Catherine Tong, Sergei Yakneen, Ellen D. Zhong, Michal
219 Zielinski, Augustin Žídek, Victor Bapst, Pushmeet Kohli, Max Jaderberg, Demis Hassabis,
220 and John M. Jumper. Accurate structure prediction of biomolecular interactions with al-
221 phafold 3. *Nature*, 630(8016):493–500, 2024. doi: 10.1038/s41586-024-07487-w. URL
222 <https://doi.org/10.1038/s41586-024-07487-w>.
- 223 N. Bennett, Joseph L. Watson, R. Ragotte, Andrew J. Borst, DéJenaé L. See, Connor Weidle, Riti
224 Biswas, Yutong Yu, Ellen L. Shrock, Russell Ault, P. J. Leung, Buwei Huang, Inna Goreshnik,
225 John Tam, Kenneth D. Carr, Benedikt Singer, Cameron Criswell, B. Wicky, Dionne K. Vafeados,
226 Mariana Garcia Sanchez, Ho Min Kim, Susana Vázquez Torres, Sidney Chan, Shirley M. Sun,
227 Timothy T Spear, Yi Sun, Keelan O’Reilly, John M. Maris, Nikolaos G. Sgourakis, Roman A.
228 Melnyk, Chang C. Liu, and David Baker. Atomically accurate de novo design of antibodies with
229 rfdiffusion. *bioRxiv*, 2024.
- 230 Chai Discovery. Chai-1: Decoding the molecular interactions of life. *bioRxiv*, 2024. doi: 10.
231 1101/2024.10.10.615955. URL [https://www.biorxiv.org/content/early/2024/
232 10/11/2024.10.10.615955](https://www.biorxiv.org/content/early/2024/10/11/2024.10.10.615955).
- 233 Felipe Leno da Silva, Andre Gonçalves, Sam Nguyen, Denis Vashchenko, R. Glatt, Thomas A.
234 Desautels, Mikel Landajuela, Daniel M. Faissol, and Brenden K. Petersen. Language model-
235 accelerated deep symbolic optimization. *Neural Computing and Applications*, 37:25601 – 25617,
236 2023.
- 237 Kerr Ding, Michael Chin, Yunlong Zhao, Wei Huang, Binh Khanh Mai, Huanan Wang, Peng Liu,
238 Yang Yang, and Yunan Luo. Machine learning-guided co-optimization of fitness and diversity
239 facilitates combinatorial library design in enzyme engineering. *Nature Communications*, 15, 2024.
- 240 Richard Evans, Michael O’Neill, Alexander Pritzel, Natasha Antropova, Andrew Senior, Tim Green,
241 Augustin Žídek, Russ Bates, Sam Blackwell, Jason Yim, Olaf Ronneberger, Sebastian Boden-
242 stein, Michal Zielinski, Alex Bridgland, Anna Potapenko, Andrew Cowie, Kathryn Tunyasuvunakool,
243 Rishub Jain, Ellen Clancy, Pushmeet Kohli, John Jumper, and Demis Hassabis. Protein complex
244 prediction with alphafold-multimer. *bioRxiv*, 2022. doi: 10.1101/2021.10.04.463034.
- 245 Christopher Frank, Ali Khoshouei, Lara Fuß, Dominik Schiwietz, Dominik Putz, Lara Weber,
246 Zhixuan Zhao, Motoyuki Hattori, Shihao Feng, Yosta de Stigter, Sergey Ovchinnikov, and Hendrik
247 Dietz. Scalable protein design using optimization in a relaxed sequence space. *Science*, 386:439 –
248 445, 2024.
- 249 Nate Gruver, S. Stanton, Nathan C Frey, Tim G. J. Rudner, I. Hotzel, J. Lafrance-Vanasse, A. Ra-
250 jpal, Kyunghyun Cho, and A. Wilson. Protein design with guided discrete diffusion. *ArXiv*,
251 abs/2305.20009, 2023.
- 252 Chak Shing Lee, Conor F. Hayes, Denis Vashchenko, and Mikel Landajuela. Reinforcement learning
253 for antibody sequence infilling. *bioRxiv*, 2025.
- 254 Luis S. Mille-Fragoso, John N. Wang, Claudia L. Driscoll, Haoyu Dai, Talal Widatalla, Xiaowei
255 Zhang, Brian L. Hie, and Xiaojing J. Gao. Efficient generation of epitope-targeted de novo
256 antibodies with germinal. *bioRxiv*, 2025.
- 257 Martin Pacesa, Lennart Nickel, Christian Schellhaas, Joseph Schmidt, Ekaterina Pyatova, Lucas
258 Kissling, Patrick Barendse, Jagrity Choudhury, Srajan Kapoor, A. Alcaraz-Serna, Yehlin Cho,
259 Kourosh H. Ghamary, Laura Vinué, B. Yachnin, Andrew M. Wollacott, Stephen Buckley, A. West-
260 phal, S. Lindhoud, S. Georgeon, Casper A. Goverde, G. Hatzopoulos, P. Gönczy, Yannick D.
261 Muller, Gerald Schwank, D. C. Swarts, A. Vecchio, Bernard L. Schneider, Sergey Ovchinnikov,
262 and Bruno E. Correia. One-shot design of functional protein binders with bindcraft. *Nature*, 646:
263 483 – 492, 2025.
- 264 Richard W. Shuai, Jeffrey A. Ruffolo, and Jeffrey J. Gray. Iglm: Infilling language modeling for
265 antibody sequence design. *Cell systems*, 2023.

270 Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, et al. De novo
271 design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023. doi:
272 10.1038/s41586-023-06415-8.
273

274 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
275 learning. *Machine Learning*, 8(3):229–256, 1992. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
276

277 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
278 Gradient surgery for multi-task learning, 2020. URL [https://arxiv.org/abs/2001.](https://arxiv.org/abs/2001.06782)
279 [06782](https://arxiv.org/abs/2001.06782).
280

281 Wensi Zhu, Aditi Shenoy, Petras Kundrotas, and Arne Elofsson. Evaluation of alphafold-multimer
282 prediction on multi-chain protein complexes. *Bioinformatics*, 39(7):btad424, 07 2023. ISSN
283 1367-4811. doi: 10.1093/bioinformatics/btad424. URL [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/btad424)
284 [bioinformatics/btad424](https://doi.org/10.1093/bioinformatics/btad424).
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

A APPENDIX

A.1 RELATED WORK

Logit representations for protein sequence design. Modern hallucination pipelines Frank et al. (2024); Pacesa et al. (2025); Mille-Fragoso et al. (2025) often optimize continuous logits $\phi \in \mathbb{R}^{N \times M}$ that induce a factorized categorical policy over sequences:

$$p_{i,j} = \frac{\exp(\phi_{i,j})}{\sum_{k=1}^M \exp(\phi_{i,k})}, \quad p(\mathbf{x} | \phi) = \prod_{i=1}^N p_{i,x_i}. \quad (1)$$

Equation (1) is the common parameterization used in logit-based sequence optimization and underlies both differentiable hallucination updates and policy-gradient updates.

Policy gradients optimize via the score-function estimator:

$$\nabla_{\phi} J = \mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [(R(\mathbf{x}) - b) \nabla_{\phi} \log(p(\mathbf{x} | \phi))]$$

enabling optimization of black-box objectives.

RL objective and policy-gradient approaches. When design quality is measured by a black-box reward $R(\mathbf{x})$, the objective and its score-function gradient can be written jointly as

$$J(\phi) = \mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [R(\mathbf{x})], \quad \nabla_{\phi} J = \mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [R(\mathbf{x}) - b \nabla_{\phi} \log p(\mathbf{x} | \phi)] \quad (2)$$

with baseline b for variance reduction. RL-style sequence optimization has been used for difficult, multi-objective fitness landscapes and diversity-constrained design Ding et al. (2024); da Silva et al. (2023); Lee et al. (2025).

Hallucination trick: moving from expectation to deterministic surrogate. Hallucination-style design commonly replaces the stochastic expected objective Eq. (2) with a differentiable objective evaluated on soft sequence representations. Let R denote a sequence-scoring functional and let p_{ϕ} be induced by logits in Eq. (1). The usual approximation is

$$J(\phi) = \mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [R(\mathbf{x})] \approx R(\mathbb{E}_{\mathbf{x} \sim p(\cdot | \phi)} [\mathbf{x}]) \approx \hat{R}(p_{\phi}), \quad \nabla_{\phi} J \approx \nabla_{\phi} \hat{R}(p_{\phi}). \quad (3)$$

where \hat{R} is differentiable with respect to soft marginals. This is accurate when distributional variance is sufficiently small or R is locally well-approximated around the mean representation, which is exactly the regime targeted by entropy reduction / sharpened logits in many hallucination pipelines Frank et al. (2024); Pacesa et al. (2025); Mille-Fragoso et al. (2025).

Why combine both. The two paradigms are complementary: Eq. (3) gives low-variance, efficient gradients for differentiable objectives, while Eq. (2) includes the score-function term that enables optimization of non-differentiable black-box rewards Ding et al. (2024); Frank et al. (2024). We refer to this hybrid as **RALO** (*RL-Augmented Logit Optimization*): it retains the hallucination backbone and injects an optional RL gradient term only in selected optimization regions, yielding a practical method that preserves baseline behavior when RL is disabled.

A.2 PSEUDO-CODE FOR RL-AUGMENTED GERMINAL STEP

A.3 POLICY GRADIENT HYPERPARAMETERS

For policy gradient optimization, we keep the RL term small ($\lambda_{\text{RL}} = 0.02$) and sparse (every 20 optimization steps), and we only start RL updates after the optimization has warmed up and the distribution is sufficiently soft. Our default is to allow a 30-gradient-step warm-up period and require distribution softness of 0.8.

A.4 RELAXED FILTER SET

For our experiments, we use a relaxed set of Germinal’s original filters.

A.5 REWARD FUNCTION WEIGHTS

Here we list the weights placed on the reward function derived from Chai-1 outputs.

We also apply reward clipping with $c = 2.0$.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Algorithm 1 RL-augmented Germinal step

Require: current logits ϕ_t , Germinal gradient tensor g_t (aux["grad"] ["seq"]), step index t , softness s_t

Require: design mode $m_t \in \{\text{logits}, \text{soft}, \text{greedy}\}$

Require: RALO hyperparameters: $\lambda_{\text{RL}} = 0.02$, $K = 10$, frequency $F = 20$, start step $t_0 = 30$, minimum softness $s_{\text{min}} = 0.8$

```

1: if RL disabled then
2:   Use Germinal update only
3:   return
4: end if
5: if  $m_t \notin \{\text{logits}, \text{soft}\}$  or  $t < t_0$  or  $s_t < s_{\text{min}}$  or  $t \bmod F \neq 0$  then
6:   Skip RL update for this step
7: else
8:   Sample  $K$  sequences  $\mathbf{x}^{(k)} \sim p(\cdot | \phi_t)$  with stochastic categorical sampling
9:   for  $k = 1, \dots, K$  do
10:    Run Chai-1 prediction for  $\mathbf{x}^{(k)}$ 
11:    Extract  $\{m_{\text{plddt}}^{(k)}, m_{\text{iptm}}^{(k)}, m_{\text{pae}}^{(k)}\}$ 
12:    Convert each metric to margin against its corresponding final-filter threshold
13:     $R^{(k)} \leftarrow 1.0 \cdot \Delta_{\text{plddt}}^{(k)} + 1.0 \cdot \Delta_{\text{iptm}}^{(k)} - 0.2 \cdot \Delta_{\text{pae}}^{(k)}$ 
14:     $R^{(k)} \leftarrow \text{clip}(R^{(k)}, -2.0, 2.0)$ 
15:   end for
16:   if  $\text{std}(\{R^{(k)}\}_{k=1}^K) \leq \epsilon$  then
17:     Skip RL update (rl_skip_no_variance=true)
18:   else
19:      $b \leftarrow \frac{1}{K} \sum_{k=1}^K R^{(k)}$ 
20:      $g_{\text{RL}} \leftarrow \frac{1}{K} \sum_{k=1}^K (R^{(k)} - b) \nabla_{\phi} \log p(\mathbf{x}^{(k)} | \phi_t)$ 
21:      $g_t \leftarrow g_t + \lambda_{\text{RL}} g_{\text{RL}}$ 
22:   end if
23: end if
24: Apply the standard Germinal optimizer step using updated  $g_t$ 

```

Parameter	Value
RL enabled	enable_rl_structural_term=true
Reward backend	Chai-1 (structure_model=chai)
Use external prediction	rl_use_structure_prediction=true
Reward components	[external_plddt, external_iptm, external_pae]
Reward weights	{external_plddt:1, external_iptm:1, external_pae:-0.2}
Reward mode	margin (relative to filter thresholds)
Reward clip	rl_reward_clip=2.0
Allow soft-mode updates	rl_apply_in_soft=true
RL coefficient	$\lambda_{\text{RL}} = 0.02$
Update frequency	rl_update_frequency=20
Samples per update	rl_num_samples=10
Sampling strategy	stochastic
Warm-up	rl_start_step=30
Min softness	rl_min_soft=0.8

Table 1: RL configuration used for the PD-L1 Germinal extension experiments.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Filter	Condition
Clashes	< 2
self-consistency RMSD (scRMSD)	< 8.0
Binder near hotspot	True
CDR3 Hotspot contacts	> 0
Percent Interface is CDR	> 40%
Interface shape complementarity	≥ 0.5
Interface hydrogen bonds	> 2
Surface hydrophobicity	≤ 0.5
Interface hydrophobicity	≥ 40
pDockQ2 Zhu et al. (2023)	> 0.15
External pLDDT	> 0.82
External ipTM	> 0.60
External pTM	> 0.75
External PAE	< 12.0

Table 2: Filter set used for our experiments.

Metric	Weight
pLDDT	+1
ipTM	+1
PAE	-0.2

Table 3: Weighting terms placed on metrics used in the reward function.