# Supplementary Information for
# MolEncoder: Towards Optimal Masked Language Modeling for Molecules

Fabian P. Krüger[1,2,3]✉, Nicklas Österbacka[1], Mikhail Kabeshov[1], Ola Engkvist[1,4], Igor Tetko[3]

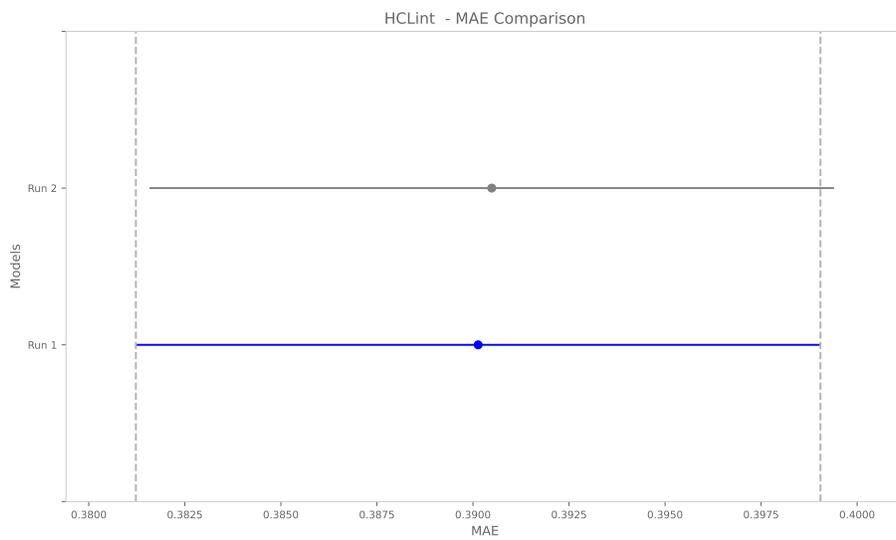[1]AstraZeneca R&D, Discovery Sciences, Molecular AI, 431 83 Mölndal, Sweden
[2]TUM School of Computation, Information and Technology, Department of Mathematics, Technical University of Munich, 80333 Munich, Germany
[3]Helmholtz Munich – German Research Center for Environmental Health (GmbH), Institute of Structural Biology, Molecular Targets and Therapeutics Center, 85764 Neuherberg, Germany
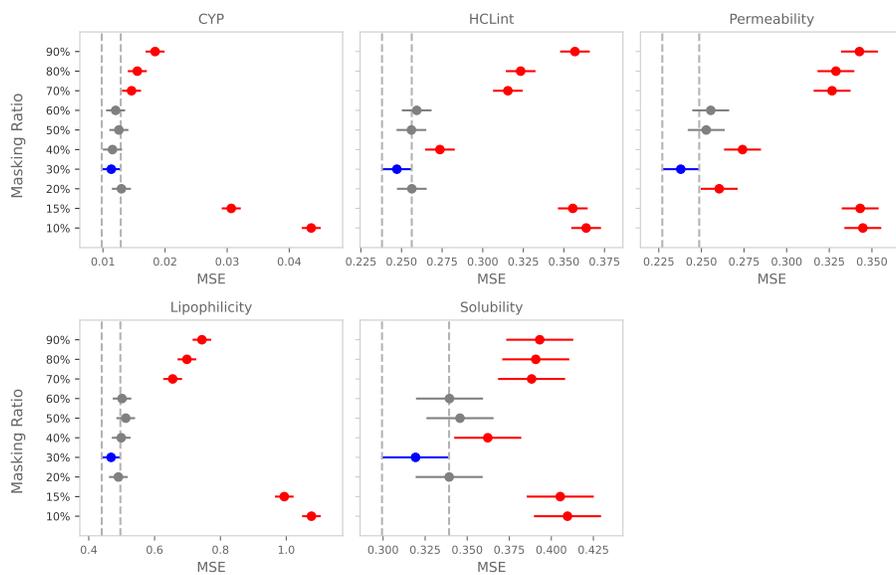[4]Department of Computer Science and Engineering, Chalmers University of Technology and University of Gothenburg, Sweden
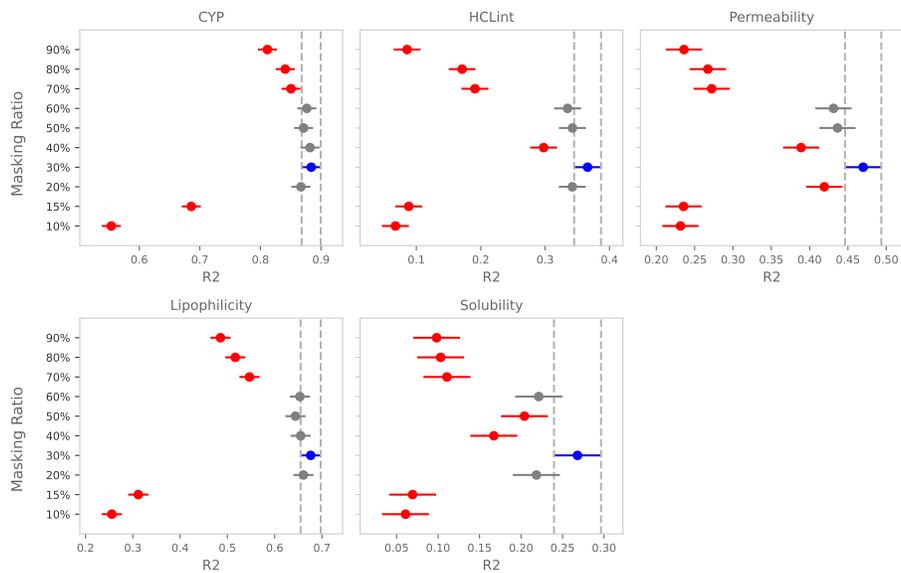
✉ fabian.krueger@tum.de
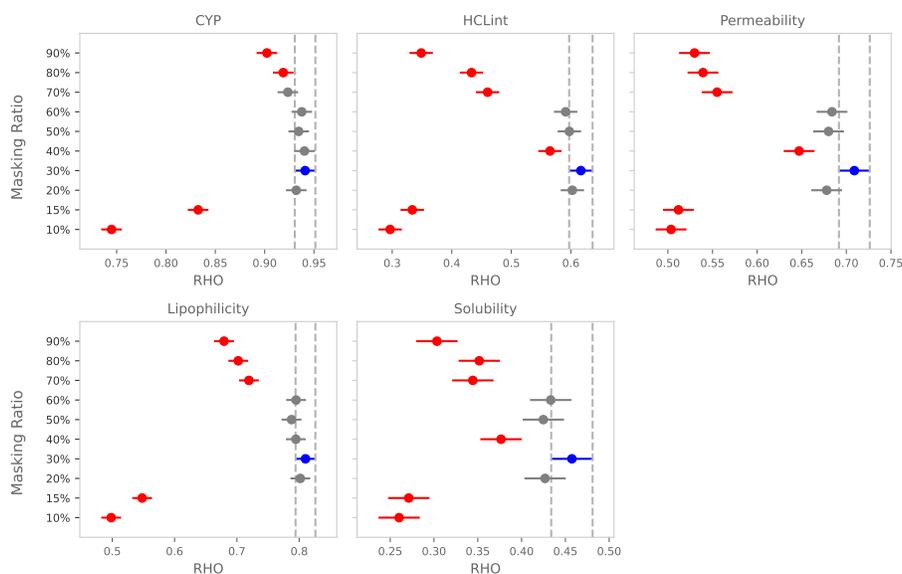
October 2025

**Supplementary Figure S1.** Comparison of two runs of the same model in our evaluation pipeline using different random seeds. The null hypothesis that both models have the same performance is not rejected, demonstrating the robustness of the evaluation pipeline to random seed variation (overlapping 95% confidence interval shown in grey).

**Supplementary Figure S2.** Same experiment as Figure 1 in the main document, but with mean squared error as an evaluation metric.



**Supplementary Figure S3.** Same experiment as Figure 1 in the main document, but with $R^2$ score as an evaluation metric.

**Supplementary Figure S4.** Same experiment as Figure 1 in the main document, but with Spearman correlation as an evaluation metric.



**Supplementary Figure S5.** Comparison of pretraining and evaluating on SMILES strings with explicit hydrogen atoms against SMILES strings with implicit hydrogen atoms (standard SMILES strings). The models that are compared have 15M parameters and are pretrained on the molecules in ChEMBL.

**Supplementary Figure S6.** Same experiment as Figure 1 in the main document, but with a model pretrained and evaluated on SMILES strings with explicit hydrogen atoms.

**HCLint - R2**

| Model Size | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.237 | 0.278 | 0.302 | 0.235 |
| 15M | 0.214 | 0.400 | 0.391 | 0.268 |
| 111M | 0.086 | 0.353 | 0.318 | 0.341 |

Pretraining Dataset Size

**Lipophilicity - R2**

| Model Size | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.371 | 0.611 | 0.628 | 0.539 |
| 15M | 0.384 | 0.692 | 0.685 | 0.632 |
| 111M | 0.171 | 0.641 | 0.635 | 0.647 |

Pretraining Dataset Size

**Permeability - R2**

| Model Size | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.335 | 0.402 | 0.399 | 0.352 |
| 15M | 0.337 | 0.484 | 0.450 | 0.399 |
| 111M | 0.246 | 0.377 | 0.400 | 0.440 |

Pretraining Dataset Size

**Solubility - R2**

| Model Size | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.188 | 0.180 | 0.219 | 0.169 |
| 15M | 0.165 | 0.276 | 0.277 | 0.161 |
| 111M | 0.081 | 0.223 | 0.223 | 0.201 |

Pretraining Dataset Size

**CYP - R2**

| Model Size | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.646 | 0.857 | 0.857 | 0.826 |
| 15M | 0.497 | 0.886 | 0.889 | 0.865 |
| 111M | 0.317 | 0.861 | 0.882 | 0.870 |

Pretraining Dataset Size

Blue = Best model
Grey = Not significantly worse
Red = Significantly worse

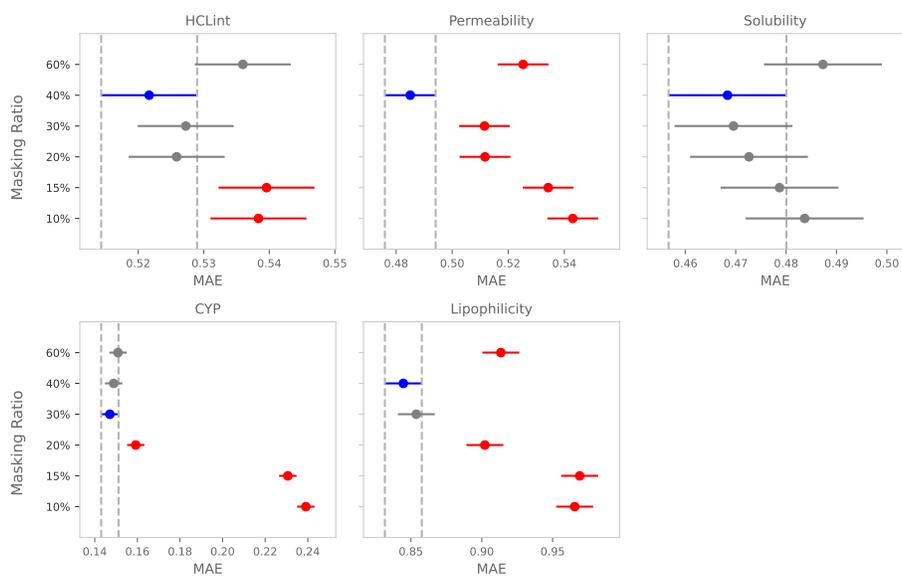**Supplementary Figure S7.** Same experiment as Figure 2 in the main document, but with $R^2$ score as an evaluation metric.

6

**HCLint - RHO**
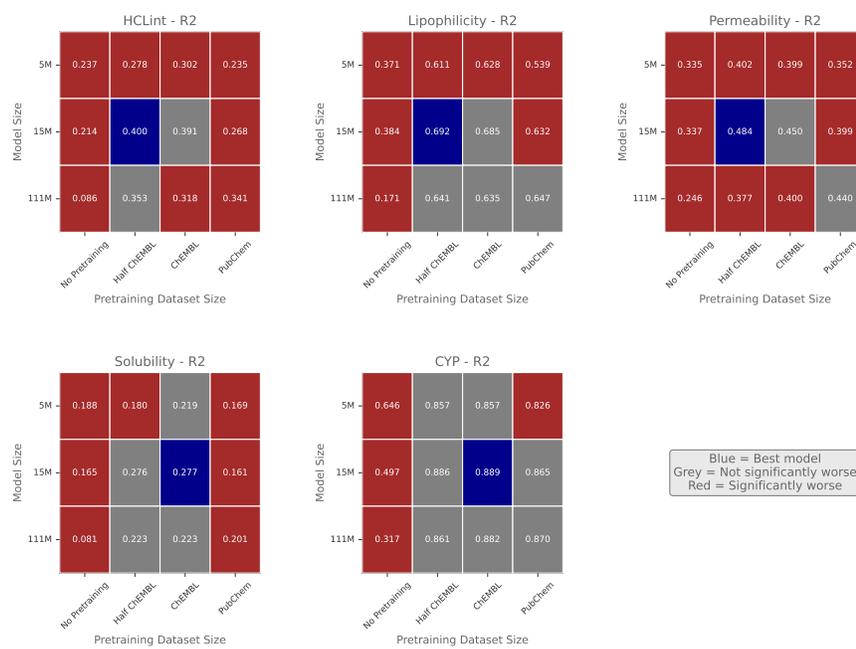
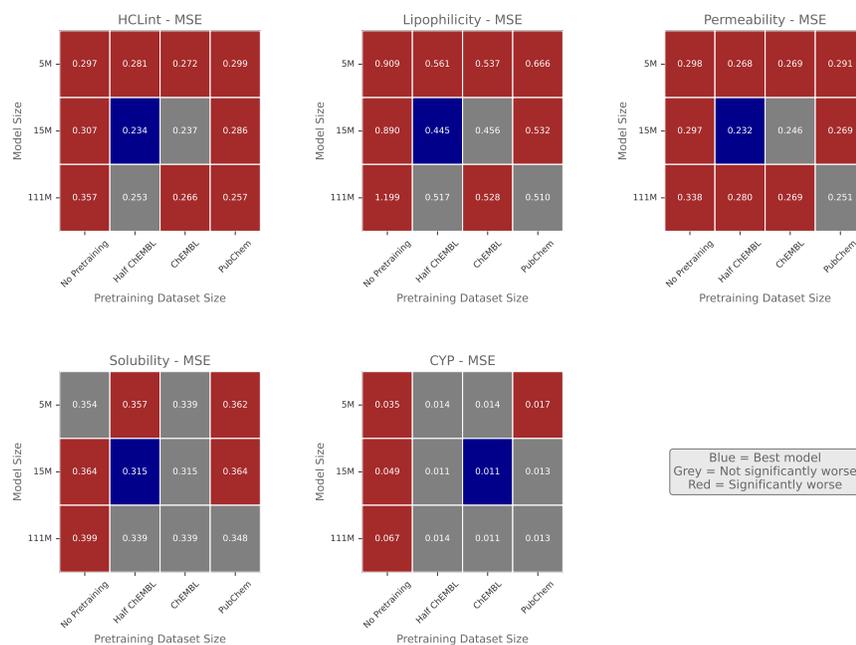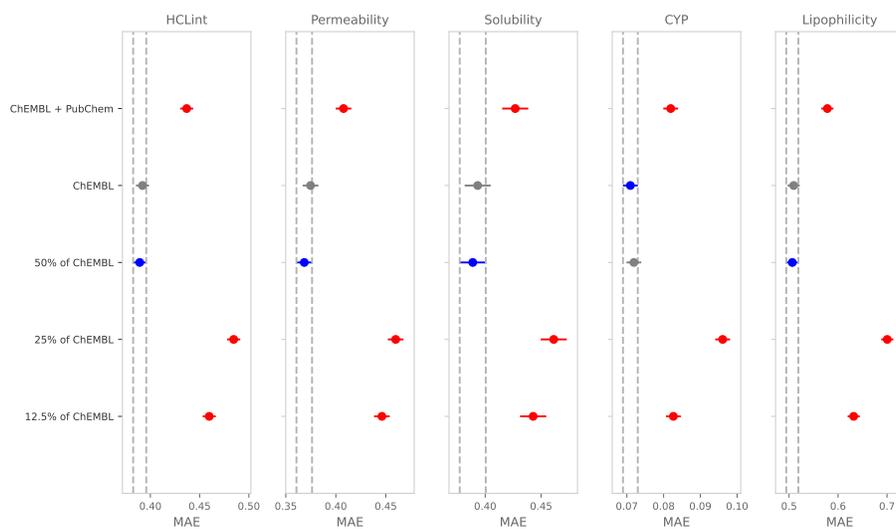|  | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.507 | 0.554 | 0.572 | 0.506 |
| 15M | 0.486 | 0.635 | 0.634 | 0.545 |
| 111M | 0.313 | 0.594 | 0.568 | 0.595 |

**Lipophilicity - RHO**

|  | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.609 | 0.768 | 0.780 | 0.715 |
| 15M | 0.605 | 0.819 | 0.814 | 0.781 |
| 111M | 0.399 | 0.785 | 0.778 | 0.790 |

**Permeability - RHO**

|  | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.600 | 0.663 | 0.660 | 0.616 |
| 15M | 0.605 | 0.716 | 0.694 | 0.662 |
| 111M | 0.507 | 0.623 | 0.648 | 0.684 |

**Solubility - RHO**

|  | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.363 | 0.382 | 0.419 | 0.377 |
| 15M | 0.348 | 0.481 | 0.471 | 0.389 |
| 111M | 0.265 | 0.434 | 0.437 | 0.426 |

**CYP - RHO**

|  | No Pretraining | Half ChEMBL | ChEMBL | PubChem |
|---|---|---|---|---|
| 5M | 0.812 | 0.927 | 0.927 | 0.910 |
| 15M | 0.706 | 0.942 | 0.944 | 0.931 |
| 111M | 0.567 | 0.927 | 0.941 | 0.934 |

Blue = Best model
Grey = Not significantly worse
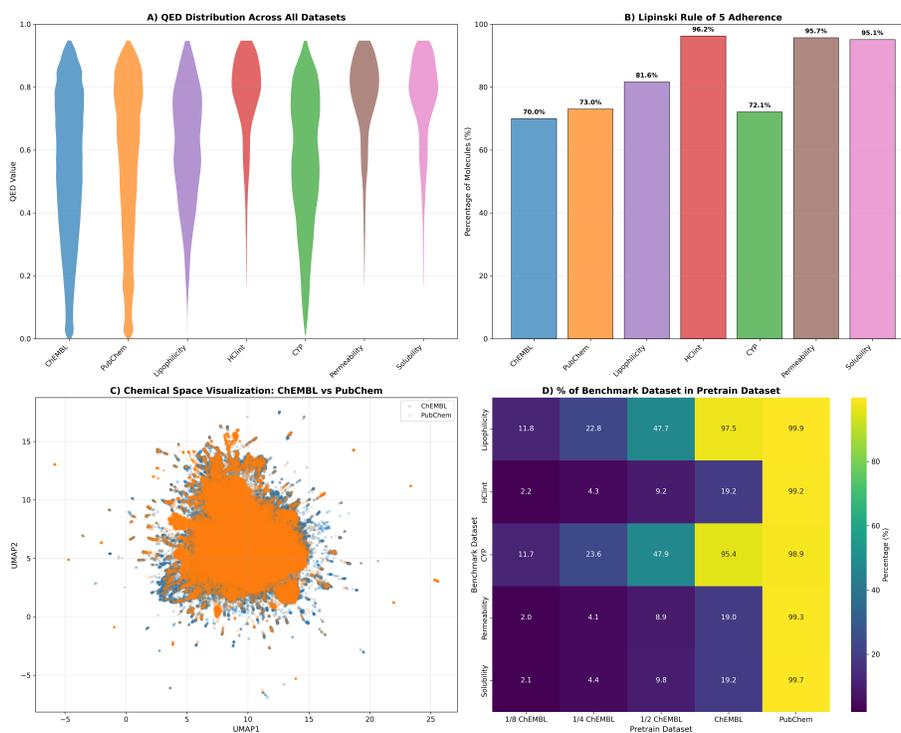Red = Significantly worse

**Supplementary Figure S8.** Same experiment as Figure 2 in the main document, but with Spearman correlation as an evaluation metric.
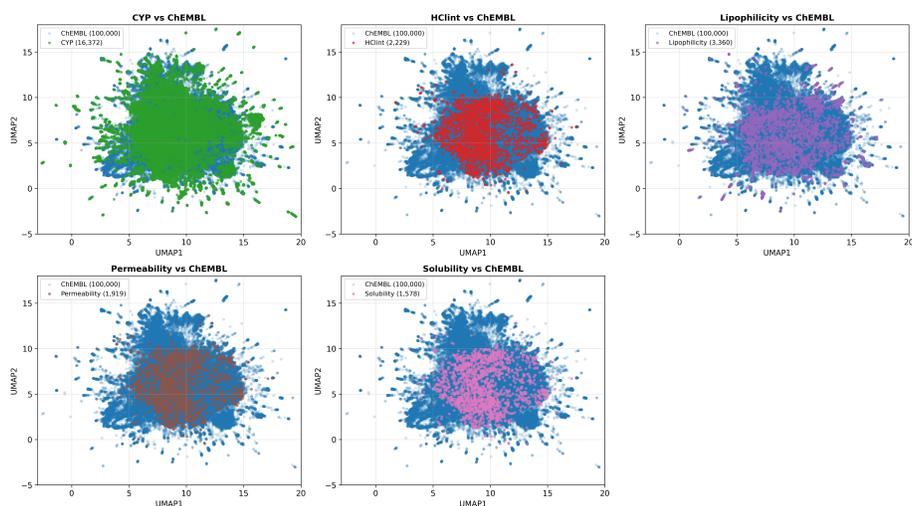
7

**Supplementary Figure S9.** Same experiment as Figure 2 in the main document, but with mean squared error as an evaluation metric.
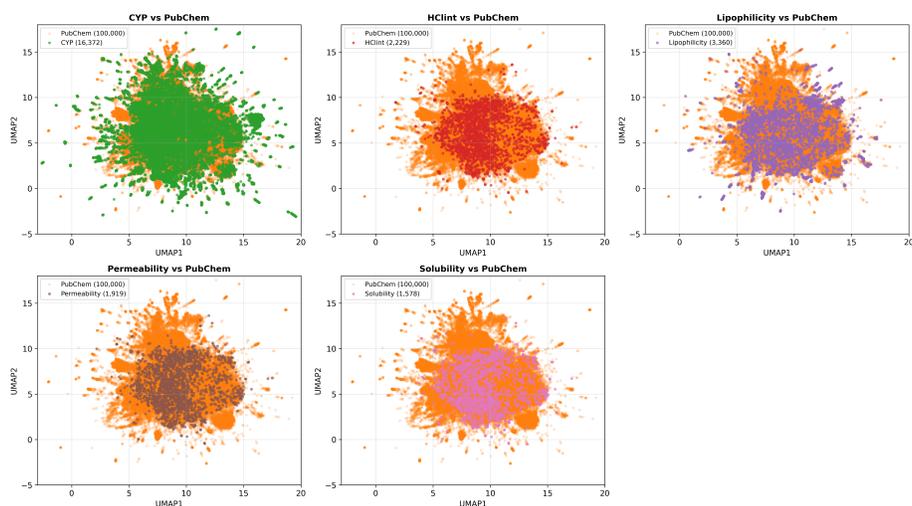


**Supplementary Figure S10.** Comparison of models pretrained on differently sized datasets using masked language modeling with a masking ratio of 30%. The models used for this comparison had 15M parameters.
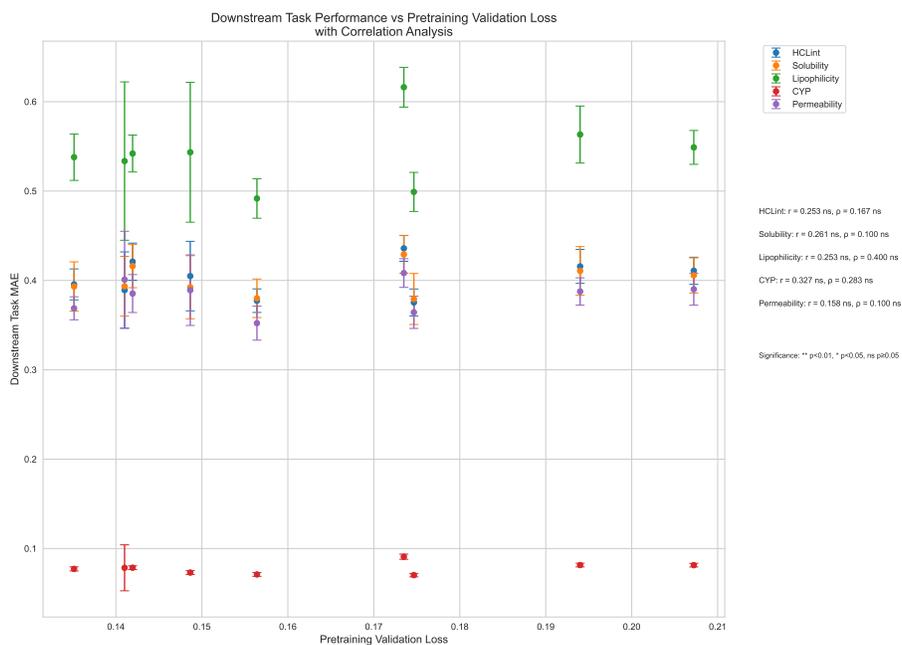
**Supplementary Figure S11.** Comparison of molecular property and chemical space distributions between ChEMBL and PubChem. (A) Quantitative estimate of drug-likeness (QED) distributions. (B) Fraction of molecules satisfying Lipinski's Rule of Five. (C) Two-dimensional UMAP projection of molecular fingerprints constructed from the first 200 principal components, explaining approximately 55% of the variance. (D) Fraction of SMILES strings from each pretraining dataset that are also present in the downstream evaluation datasets. Both datasets were uniformly subsampled to 100,000 molecules for computational efficiency for the subplots A-C. Subplot D contains the full datasets.
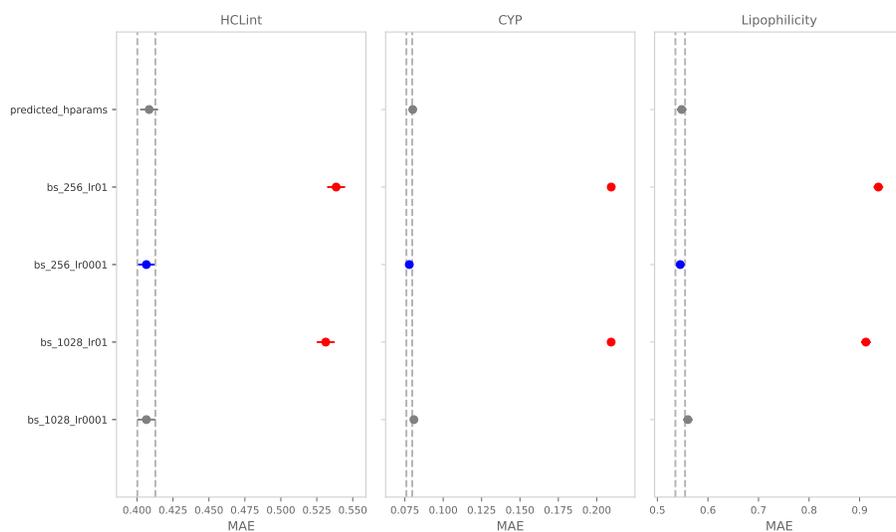
**Supplementary Figure S12.** Comparison of chemical space between ChEMBL and the downstream evaluation datasets. Each point represents a molecule embedded in a two-dimensional UMAP projection derived from the first 200 principal components of molecular fingerprints (explaining 55% of the variance). The same UMAP embedding used for the ChEMBL–PubChem comparison was applied here to ensure a consistent latent space.

**Supplementary Figure S13.** Comparison of chemical space between PubChem and the downstream evaluation datasets. Each point represents a molecule embedded in a two-dimensional UMAP projection derived from the first 200 principal components of molecular fingerprints (explaining 55% of the variance). The same UMAP embedding used for the ChEMBL–PubChem comparison was applied here to ensure a consistent latent space.
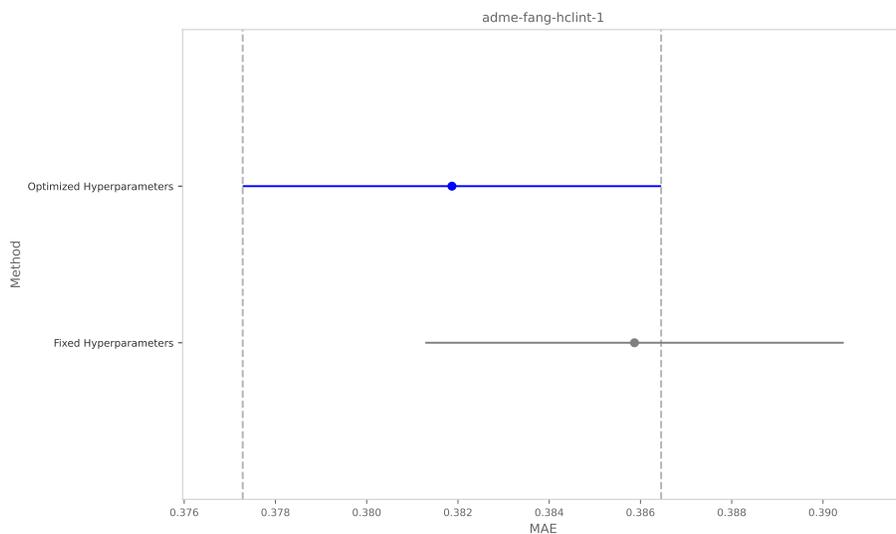
**Supplementary Figure S14.** Downstream performance (measured as MAE) plotted against pretraining performance (cross-entropy loss). Pearson and Spearman correlation coefficients for each evaluation dataset are shown next to the figure. No significant correlation was observed for any of the evaluation datasets.

**Supplementary Figure S15.** Comparison of models pretrained with different hyperparameters (batch size and learning rate). The predicted hyperparameters correspond to those used in the study and do not differ significantly from the best values identified through manual hyperparameter tuning.



**Supplementary Figure S16.** Comparison of fine-tuning the pretrained 15M-parameter model using our fixed fine-tuning hyperparameters versus Bayesian hyperparameter optimization for fine-tuning on the specific dataset. The resulting model performances do not differ significantly.

**Supplementary Table S1.** Mean Absolute Error (MAE) and benchmark-specific ranks for *MolEncoder*. The rank reflects the leaderboard at the 03.07.2025. Archived versions of the leaderboard are available at
`https://web.archive.org/web/20250703093745/https://polarishub.io/benchmarks/polaris/adme-fang-perm-1`,
`https://web.archive.org/web/20250703094412/https://polarishub.io/benchmarks/polaris/adme-fang-hclint-1`,
`https://web.archive.org/web/20250703094653/https://polarishub.io/benchmarks/polaris/adme-fang-solu-1`,
`https://web.archive.org/web/20250703094922/https://polarishub.io/benchmarks/novartis/adme-novartis-cyp3a4-reg`,
`https://web.archive.org/web/20250703095138/https://polarishub.io/benchmarks/tdcommons/lipophilicity-astrazeneca`.

| Benchmark | MAE | Rank (out of N) |
|---|---|---|
| Permeability | 0.305 | 3 / 13 |
| HCLint | 0.337 | 4 / 12 |
| Solubility | 0.380 | 13 / 43 |
| CYP | 0.198 | 2 / 2 |
| Lipophilicity | 0.497 | 4 / 8 |

**Supplementary Table S2.** Model configuration details for each model variant.

| Model | Layers | Hidden Size | Heads | Intermediate Size |
|---|---|---|---|---|
| 5M | 8 | 256 | 4 | 384 |
| 15M | 12 | 384 | 6 | 576 |
| 111M | 22 | 768 | 12 | 1152 |

**Supplementary Table S3.** Pretraining configuration and optimizer settings. Learning rate and batch size formulas are taken from Li et al. 2025.

| Parameter | Value / Description |
|---|---|
| Optimizer | Schedule-free AdamW |
| Learning rate | $\eta(N, D) = 1.79\, N^{-0.713}\, D^{0.307}$ |
| | $N$: number of non-embedding parameters, |
| | $D$: dataset size in tokens |
| Batch size (in tokens) | $B(D) = 0.58\, D^{0.571}$ |
| | $D$: dataset size in tokens |
| Weight decay | 0.00001 |
| Adam $\beta_1$ | 0.9 |
| Adam $\beta_2$ | 0.999 |
| Warmup steps | 1000 |
| Mixed precision | Enabled |
| Model compilation | `torch.compile` with Inductor backend |
| Metric for best model | Cross-entropy loss on held-out test set |
| Dataloader workers | 32 |
| Pin memory | True |

**Supplementary Table S4.** Hyperparameter combinations used in the pretraining ablation study. The underlined configuration is the one predicted by the scaling rules of Li et al. 2025.

| Learning Rate | Batch Size (samples) |
|---|---|
| <u>0.004255</u> | <u>512</u> |
| 0.001 | 256 |
| 0.001 | 1024 |
| 0.1 | 256 |
| 0.1 | 1024 |

**Supplementary Table S5.** Hyperparameter search space for finetuning optimization using TPE.

| Hyperparameter | Search Range |
|---|---|
| Learning rate | 1e-4 to 1e-3 (log scale) |
| Weight decay | 1e-6 to 1e-2 (log scale) |
| Warmup steps | 0 to 100 |
| Batch size | 32, 64, 128, 256, 512 |

# References

Houyi Li, Wenzhen Zheng, Qiufeng Wang, Hanshan Zhang, Zili Wang, Shijie Xuyang, Yuantao Fan, Shuigeng Zhou, Xiangyu Zhang, and Daxin Jiang. Predictable scale: Part i–optimal hyperparameter scaling law in large language model pretraining. arXiv preprint arXiv:2503.04715, 2025.