A MATRIX VARIATIONAL AUTO-ENCODER FOR VARIANT EFFECT PREDICTION IN PHARMACOGENES

Anonymous authors

Paper under double-blind review

ABSTRACT

Variant effect predictors (VEPs) are designed to predict the impact of protein variants on cellular function, traditionally using data from multiple sequence alignments (MSAs). This assumes that natural variants are fit, a premise challenged by pharmacogenomics, where some pharmacogenes have low evolutionary pressure. In this context, deep mutational scanning (DMS) datasets are of particular interest since they provide quantitative fitness scores for variants. In this work, we propose a transformer-based matrix variational auto-encoder architecture and evaluate its performances on 33 DMS datasets corresponding to 26 drug target and absorptiondistribution-metabolism-excretion (ADME) proteins available in the ProteinGym benchmark. Our model trained on MSAs (matVAE-MSA) outperforms a model similar to the widely used VEPs in pharmacogenomics, and sets a new zero-shot prediction benchmark for 2 proteins related to the Noonan syndrome. We compare matVAE-MSA with matENC-DMS, a model with similar capacity, but trained on DMS data in a 5-fold supervised cross-validation framework. matENC-DMS outperforms matVAE-MSA for 15 out of 33 DMS datasets, including all ADME, and certain drug target proteins. Although our models do not outperform the best baseline models, our results help shed new light on the role of evolutionary pressure for the validity of the premise of VEP design. In turn motivating the development of DMS datasets to improve VEPs on pharmacogene-related proteins.

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

033 Variant effect predictors (VEPs) are mathematical models aiming at predicting the effect of one or 034 multiple variants in a sequence of amino-acids (AAs). The effect of a protein variant is typically defined as a loss or gain of function of a cell carrying the variant, compared to a cell carrying a 035 wild-type (WT) protein without variant. The accurate prediction of variant effect has many promising applications for personalized medicine, particularly in the field of pharmacogenomics, where 037 variants on drug targets or Absorption-Distribution-Metabolism-Excretion (ADME) proteins are of particular interest (Huang et al., 2016). In this context, VEPs can be used to assess individual patient response to chemotherapeutic treatments from their genetic background, thus eliminating the need 040 for multiple attempts at treatments. The most effective VEPs have been designed using data from 041 multiple sequence alignments (MSAs) and based on the conservation assumption: fit variants were 042 selected out by nature and thus, learning a distribution over variants found in nature implicitly cap-043 tures the biochemical constraints that characterize fit variants. New sequencing techniques combined 044 with machine learning could lead to significant advances in variant effect prediction, by providing quantitative data in region of the protein sequence space unexplored in existing MSA datasets. Deep mutational scanning (DMS) has recently emerged as a way to yield large-scale datasets of protein 046 quantitative fitness scores (Fowler & Fields, 2014). The fitness scores can also be obtained with 047 different selection assays, allowing to quantify various effects, e.g. effect on phenotype or effect on 048 structure. DMS thus allows to challenge the conservation assumption of VEPs design from MSAs. This is in turn is of particular importance in pharmacogenomics, since pharmacogenes are generally under low evolutionary pressure (Zhou et al., 2022; Ingelman-Sundberg et al., 2018). 051

In this article, we design a VEP and use it to evaluate the validity of the conservation assumption for
 pharmacogene-related proteins. Our architecture exploits the structure of variational auto-encoders
 (VAEs) and allows models of similar capacity and designs to be trained on both MSA and DMS

data. We also exploit a transformer architecture in order to improve upon existing VAE-based VEPs.
 Both VAE and transformers are key components of the best performing models in the ProteinGym
 benchmark (Notin et al., 2023a). We experiment with a VAE-based model exploiting multimodal
 priors, and we derive a matrix encoding scheme inspired from linear matrix decomposition to replace
 the input flattening operation found for instance in DeepSequence.

060 061

062

063

064

065

066

067

068

069

070

071

073

074

082 083 084 **Contributions** Our contributions are summarized as follows:

- 1. We design protein specific models combining a VAE and a transformer for variant effect prediction. We study their zero-shot prediction performances on 33 deep mutational scanning (DMS) datasets of drug related and ADME proteins available in the ProteinGym benchmark.
- 2. We study the impact on performances of using expressive latent prior distributions when the models are trained on MSA data available in ProteinGym. We experiment with standard mixture of Gaussian (MOG) and VampPrior.
- 3. We adapt our models to directly predict labels from DMS datasets using a prediction head from the latent space, thus preserving our model capacity. In light of the comparison in performances of the models trained unsupervised on MSA and supervised DMS label data, we discuss the extent of the validity of the conservation assumption.

1.1 Related works

Zero-shot predictors VEPs exploiting site-independent position-wise frequencies of AAs in MSAs remain the methods of choice in pharmacology, e.g. SIFT or Polyphen-2 (Ng & Henikoff, 2003; Adzhubei et al., 2010; Durbin, 1998). However, other models can achieve much better zeroshot prediction performances on at least one pharmacogene-related protein DMS dataset (Details in Table A.6), according to the recent ProteinGym variant effect prediction benchmark (Notin et al., 2023a). Many of these models compute the functional cellular effect of a variant v compared to a wild-type sequence $\mathbf{x}^{(wt)}$, via the log-likelihood ratio:

$$\hat{y} = \ln \frac{p(\underline{\mathbf{x}}^{(v)})}{p(\mathbf{x}^{(wt)})},\tag{1}$$

085 where p(.) is a generative probability density chosen to maximize $p(\mathbf{x}^{(v)})$, for sequences $\mathbf{x}^{(v)}$ from the MSA. For instance, the Evolutionary Scale modeling (ESM) approaches (Rives et al., 2021; 087 Lin et al., 2023), rely solely on a transformer-based protein language model (PLM) for modeling the distribution over sequences in MSAs. Tranception (Notin et al., 2022a) additionally integrates predictions using position-wise frequencies of AAs in MSAs. TranceptEVE (Notin et al., 2022b) combines the Tranception model with a VAE-based model (Frazer et al., 2021) for AA sequence 091 modeling. Other methods such as Masked Inverse Folding (MIF) (Yang et al., 2023) learn to predict protein sequences from a given structure. VESPA (Marquet et al., 2022) combines protein sequence 092 embedding from PLMs with known bio-mechanical properties of AAs to predict variant effect with a linear regression model. Other model do not rely on the ratio in (1) to compute variant effect. MSA 094 Transformer (Rao et al., 2021) is based on ESM and uses axial attention to optimize a masking loss 095 over an entire MSA, rather than on individual sequences. It learns a representation of Hamming 096 distances in the MSA and the hamming distance to WT sequence is used as a proxy for variant 097 effect. GEMME (Laine et al., 2019) predicts variant effect via the distance to WT sequence in an 098 evolutionary tree. This approach shows very good performances and has several order of magni-099 tude fewer parameters than transformer-based approaches. DeepSequence (Riesselman et al., 2018) 100 introduces a VAE and approximates the distribution of input data \underline{x} (Eq. (1)) with the variational 101 evidence lower bound.

102

Supervised learning predictors Recently, several models combining DMS and MSA datasets
 have been proposed Hsu et al. (2022). The general idea is to combine sequence embeddings, e.g.
 sequence one-hot encoding, with evolutionary fitness scores from pretrained models such as ESM or
 DeepSequence. ProteinNPT is a conditional pseudo-generative model designed for exploiting DMS
 data, jointly with MSA data in a semi-supervised framework (Notin et al., 2023b). In addition to
 their novel architecture, the authors introduce several baselines consisting in exploiting prediction

scores from zero-shot prediction models pretrained on MSA, including DeepSequence and MSA
 Transformer. SPIRED is a recent framework able to predict fitness scores as well as protein structure
 (Chen et al., 2024). A pretrained ESM model is used for sequence embedding, and graph attention
 networks and multilayer perceptron are trained using DMS data in a supervised framework.

112

128 129

130 131

132

133

134

135

136 137

138 139 140

141 142

143 144 145

113 Multi-modal prior distributions for VAEs VAEs, e.g. DeepSequence, assume that the input data $\mathbf{x} \in \{0,1\}^{L \times d}$ are generated from a latent variable of a D-dimensional vector space: $\mathbf{z} \in \mathbb{R}^{D}$. The 114 latent variable is assumed drawn from a Gaussian prior $p(\mathbf{z})$, and the generative process is modeled 115 116 with a distribution $p_{\theta}(\mathbf{x}|\mathbf{z})$. The explicit modeling of the latent variable \mathbf{z} is an interesting feature of VAEs, because it allows to put a formal prior distribution on the latent space. The other mentioned 117 models do not impose such structure, although interestingly the learnt representations in ESM was 118 shown to correlate with known bio-mechanical properties of AAs (Rives et al., 2021). Multimodal 119 mixture of Gaussian (MOG) priors have been proposed as latent prior distributions for unsuper-120 vised clustering tasks (Dilokthanakul et al., 2017) with VAE. The authors used trainable mean and 121 covariances in latent space and showed through data sampling that the learnt mixture components 122 corresponded to meaningful characteristics of the input data. This was shown to have potential im-123 plications for model interpretability in biological contexts (Varolgünes et al., 2020). Further, the 124 VampPrior has been designed so that the statistics of the mixture components explicitly depend on 125 input space prototypes Tomczak & Welling (2018). This provides meaningful variables to probe for interpretability rather than using a sampling scheme. To the best of the authors knowledge, current 126 methods employing VAEs for VEP have only been designed with unimodal prior distributions. 127

2 Methods

A detailed description of the matVAE-MSA architecture is provided in section 2.1. A reduction of the architecture with similar capacity and that can be trained on DMS data: matENC-DMS, is proposed in section 2.2. The datasets that are used to train and evaluate the models are introduced in section A.1.



Figure 1: Model architecture of matVAE-MSA. DwFC: Dimensionwise fully connected layer. CE: Cross-entropy loss. FCB: Fully connected bottleneck.

147 148 149

150 151

160

146

2.1 MODEL DESCRIPTION

2.1.1 MATRIX DECOMPOSITION AND ENCODING

Matrix decomposition is a set of methods in linear algebra consisting in decomposing a matrix $\underline{\mathbf{x}} \in \mathbb{R}^{L \times d}$ (here d < L), into two (or more) matrices with interesting structure, e.g. unitary, triangular, diagonal. For instance, the QR decomposition of $\underline{\mathbf{x}} \in \mathbb{R}^{L \times d}$ is of the form $\underline{\mathbf{x}} = \mathbf{Q}[\mathbf{R}^T, \mathbf{0}^T]^T$, where $\mathbf{Q} \in \mathbb{R}^{L \times L}$ is unitary, $\mathbf{R} \in \mathbb{R}^{d \times d}$ is upper triangular and $\mathbf{0}$ is a $(L - d \times d)$ -dimensional matrix of zeros. We call matrix encoding the use of the low dimension factor, here the \mathbf{R} matrix, as a compressed representation of the input $\underline{\mathbf{x}}$. For QR decomposition, the matrix encoding can be formulated:

$$[\underline{\mathbf{s}}^T, \mathbf{0}^T]^T = \mathbf{W} \underline{\mathbf{x}} \in \mathbb{R}^{L \times d},$$
(2)

where $\mathbf{W} = \mathbf{Q}^{-1} \in \mathbb{R}^{L \times L}$ is a linear transform and $\underline{\mathbf{s}} = \mathbf{R}$. Linear decomposition methods are often related to singular value decomposition, for instance the diagonal elements of \mathbf{R} in the QR

decomposition are the singular values of \underline{x} . However, for one-hot encoded sequences of AAs, the singular values are counts of each AA in the sequence, and the singular vectors are a permutation of \mathbb{I}_d determined by the ordering of the counts of the AAs. In other words, the encoding produced with such linear methods does not account for the global or relative position of AAs in the sequence, and in particular, randomly permuting the rows of \underline{x} leads to the same encoding. To ensure that the model is flexible enough to learn a useful encoding, we propose to learn a representation of \underline{x} with a transformer, prior to reducing the first dimension to a fixed H < L with a trainable linear transform. This is formulated as follows:

170 171

172

179 180

187

193

200

214

215

$$\underline{\mathbf{x}}' = \operatorname{Transformer}(\underline{\mathbf{x}}) \in \mathbb{R}^{L \times d},$$

$$\underline{\mathbf{s}} = \mathsf{DwFC}(\underline{\mathbf{x}}') \in \mathbb{R}^{H \times d}$$

where Transformer(.) and DwFC(.) are specified in the paragraphs below.

Transformer Transformers are effective sequence models that can transfer information between any two positions within a sequence. The model we use is similar to the multi-layer encoding transformer in (Vaswani et al., 2017). Individual layers encode an input sequence $\underline{\mathbf{x}} \in \mathbb{R}^{L \times d}$ into a sequence $\underline{\mathbf{x}}' \in \mathbb{R}^{L \times d}$ as follows:

$$\underline{\mathbf{x}}_{1} = \operatorname{Norm}(\tau_{1}\underline{\mathbf{x}} + (1 - \tau_{1})\operatorname{Attn}(\underline{\mathbf{x}})), \\ \underline{\mathbf{x}}' = \operatorname{Norm}(\tau'\underline{\mathbf{x}}_{1} + (1 - \tau')\operatorname{FC}(\underline{\mathbf{x}}_{1})),$$
(3)

where Norm(.) is a Layer Normalization (Ba et al., 2016), FC(.) is a fully connected (FC) network with ReLU activations, and Attn(.) is the masked scaled dot product attention from (Vaswani et al., 2017). We use trainable $\tau_1, \tau' \in [0, 1]^2$ to control the gradient flow in the gated skip connections (He et al., 2016). Note that we do not use positional encoding since (Rives et al., 2021) showed that PLMs did not necessarily benefit from it. Instead, structure information is encoded in the mask of the attention layer (See section A.2).

Dimension-wise FC (DwFC) We call DwFC the FC linear layer inspired from Eq. 2. This layer is similar to a flattening followed by a linear transform with bias, but requires less parameters since the same linear transform is used across dimensions. This operation replaces the direct flattening of the input in DeepSequence. $\underline{\mathbf{x}}' \in \mathbb{R}^{L \times d}$ is encoded in a protein length independent representation $\underline{\mathbf{s}} \in \mathbb{R}^{H \times d}$ as follows:

$$\underline{\mathbf{s}} = \mathbf{U}\underline{\mathbf{x}}' + \mathbf{b},\tag{4}$$

where $\mathbf{U} \in \mathbb{R}^{H \times L}$ and $\mathbf{b} \in \mathbb{R}^{H}$ are trainable weight and bias parameters.

196 2.1.2 FULLY CONNECTED BOTTLENECK (FCB)

The FCB is similar to a classical VAE (Kingma & Welling, 2014). $\underline{s} \in \mathbb{R}^{H \times d}$ is flattened and then encoded into a latent representation of fixed dimension:

$$\mathbf{h} = \mathrm{FC}(\mathrm{Vec}(\underline{\mathbf{s}})) \in \mathbb{R}^D, \tag{5}$$

where Vec(.) denotes the flattening operation. **h** is then used to compute the statistics of the latent vector $\mathbf{z} \in \mathbb{R}^{D}$:

$$q_{\phi}(\mathbf{z}|\underline{\mathbf{x}}) = \mathcal{N}\left(\mathbf{z}; f_{\mu}(\mathbf{h}), f_{\sigma}(\mathbf{h})\right), \tag{6}$$

where **h** is a function of $\underline{\mathbf{x}}$ and f_{μ} , respectively f_{σ} , is a 1-layer fully connected network returning a mean vector, respectively a diagonal covariance matrix. For training, we introduce robustness by drawing the latent vector $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\underline{\mathbf{x}})$ using the reparameterization trick. At test time, we use $\mathbf{z} = f_{\mu}(\mathbf{h})$ to reduce stochasticity. The latent distribution $q_{\phi}(\mathbf{z}|\underline{\mathbf{x}})$ is learnt to be close to a prior distribution $p(\mathbf{z})$ with respect to the Kullback-Leibler divergence (KLD). The latent vector $\mathbf{z} \in \mathbb{R}^D$ is then used as input to a decoder network that aims to reconstruct the input $\underline{\mathbf{s}}$. The output to the decoder, and thus of the FCB, is denoted $\underline{\mathbf{\hat{s}}} \in \mathbb{R}^{H \times d}$.

Mixture of Gaussian (MOG) Prior To extend the work carried out in DeepSequence, we choose
 prior distributions formulated as a MOG:

 $p(\mathbf{z}) = \frac{1}{M} \sum_{k=1}^{M} p_k(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$ (7)

where for $k = 1, \dots, M$, $p_k(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ are multivariate Normal distributions with diagonal covariance $\boldsymbol{\Sigma}_k \in \mathbb{R}^{D \times D}$ and mean $\boldsymbol{\mu}_k \in \mathbb{R}^D$. As opposed to Gaussian distribution, MOG are multimodal and add more structure in latent space. This in turn can lead to a more expressive generative model, with a latent space able to important differences in input space in different modes, and fully use individual modes to encode subtle differences.

VampPrior The VAMP prior (Tomczak & Welling, 2018) is a special case of MOG with means and covariance that are functions of trainable prototypes in input space. That is, for k = 1, ..., M, prototypes $\underline{\mathbf{u}}_k \in \{0, 1\}^{L \times d}$ are used to compute

 $\mathbf{h}_{k} = \text{FC} \left(\text{DwFC} \left(\text{Transformer} \left(\underline{\mathbf{u}}_{k} \right) \right) \right),$

where DwFC(.) and Transformer(.) are defined in 2.1.1, and FC(.) is the fully connected encoding defined in (5). The mean and covariance of each mixture components are then computed with f_{μ} and f_{σ} , similarly to 6. The prior distribution is finally expressed as:

$$p(\mathbf{z}) = \frac{1}{M} \sum_{k=1}^{M} \mathcal{N}(\mathbf{z}; f_{\mu}(\mathbf{h}_k), f_{\sigma}(\mathbf{h}_k)).$$

where for k = 1, ..., M, \mathbf{h}_k depends upon the k-th trainable prototype $\underline{\mathbf{u}}_k$. In addition to the potential benefits of multi-modal priors mentioned before, the VampPrior could provide a way to interpret the modes of the prior distribution, in light of prototypes in input space.

2.1.3 DECODING

221

225 226

227

228

233

234

235 236

237

250 251

253

254

255

256

257 258 259

264

265

The decoding process denoted $p_{\theta}(\mathbf{x}|\mathbf{z})$, is symmetrical to the encoder, i.e. consists of a decoding FC layer, and a dimension wise FC layer followed by a transformer. One important difference is that the decoding transformer includes a temperature softmax output operation to ensure that the rows of the reconstructed $\hat{\mathbf{x}} \in \mathbb{R}^{L \times d}$ define proper discrete distributions.

242 2.1.4 Loss function

Our loss function is directly derived from the negative evidence lower bound (ELBO) used in VAEs (Kingma & Welling, 2014). In the ELBO, one term corresponds to an expected reconstruction loss (defined in section A.3). The other term is the KLD between the approximated posterior distribution q_{ϕ} , a Gaussian, and the prior p, a mixture of Gaussian distributions. The KLD between q_{ϕ} and phas no closed form expression and is therefore approximated with an upper bound (Durrieu et al., 2012):

$$D_{KL}\left(q \mid \mid \sum_{k=1}^{M} w_k p_k\right) \le -\ln \sum_{k=1}^{M} w_k e^{-D_{KL}(q \mid \mid p_k)},\tag{8}$$

where q, p_1, \ldots, p_M are distributions such that q is absolutely continuous with respect to p_1, \ldots, p_M , for $k = 1, \ldots, M$ $w_k \ge 0$ and $\sum_{k=1}^M w_k = 1$. A proof of the inequality can be found in (Rodríguez Gálvez, 2024, Appendix 6.B). Here, $q = q_{\phi}(\mathbf{z}|\mathbf{x})$, for $k = 1, \ldots, M$ $w_k = \frac{1}{M}$ and p_k are mixture components defined in (7). The negative ELBO loss function is finally written using (8) as follows:

$$l(\underline{\mathbf{x}};\theta,\phi) = -\left(\ln\frac{1}{M}\sum_{k=1}^{M}e^{-D_{KL}(q_{\phi}||p_{k})} + \mathbb{E}_{q_{\phi}(\mathbf{z}|\underline{\mathbf{x}})}\left[\ln p_{\theta}(\underline{\mathbf{x}}|\mathbf{z})\right]\right).$$
(9)

The complete encoding/decoding structure of the model depicted in Fig. 1 is trained on MSA data to minimize (9). At test time, the ELBO in (9) is used as an approximation of the log-evidence to compute the log-likelihood ratio in (1).

2.2 REDUCTION OF THE MODEL FOR DMS DATA

Our reduced model uses the encoding part of matVAE-MSA, and replaces the decoding part with a FC network prediction head to predict the quantitative DMS score $y \in \mathbb{R}$:

268
269
$$\mathbf{h} = FC(Vec(DwFC(Transformer(\underline{\mathbf{x}})))) \in \mathbb{R}^{D},$$

$$\hat{y} = FC(\mathbf{h}) \in \mathbb{R},$$

where $h \in \mathbb{R}^{D}$ is similar to (5). The reduced model depicted in Fig. 2 is referred to as "matENC-DMS". matENC-DMS has a capacity very close to that of matVAE-MSA since the encoder is identical in the two models, and the decoder of matVAE-MSA is symmetrical to the encoder and only aims at reconstructing the input. We argue that this is an ideal setup to test the conservation assumption often used when designing VEPs: unfit variants were selected out by nature and thus, learning a distribution over these sequences implicitly captures the biochemical constraints that characterize fit variants.



Figure 2: Model architecture of matENC-DMS. DwFC: Dimensionwise fully connected layer. CE: Cross-entropy loss. FC: Fully connected.

3 EXPERIMENTS

277

278

279

281

283 284

287

288

289 290 291

292 293

295

296

297 298

299

In this section we provide details on our experimental design choices. Our model comparison setup with the choice of baselines and performance metrics is explained in 3.1. The model architecture hyper-parameters are discussed in A.2, and the hyper-parameters used for training are detailed in A.3.

3.1 MODEL COMPARISON

300 **Performance metrics** The performances of our protein specific models are measured and reported 301 on corresponding individual DMS datasets, both in terms of Spearman's Rank Correlation coef-302 ficient (SpearmanR) and area under receiver operating characteristic (AUROC). The SpearmanR measures the correlation between the ranks of the predicted scores and the ranks of the target scores. 303 Additionally, the target score is binarized in order to compute the AUROC. To allow for a mean-304 ingful comparison of the AUROC scores, we use the binarization threshold used in ProteinGym. In 305 brief, given a DMS dataset, a threshold on the scores is selected manually between modes in case the 306 distribution of scores is bimodal, and as the median in case the distribution of scores is unimodal. In 307 the rest of the paper we will primarily compare spearmanR performances. 308

For matENC-DMS, we train our model on DMS data with supervised learning in a 5-fold crossvalidation framework. This first ensures that no variant/label pair is used for both training and testing. Secondly, this ensures that the evaluation framework is comparable to that of matVAE-MSA, with all variants in the DMS dataset used exactly once for validation. When a protein has multiple DMS datasets, both datasets were split in 5 folds and the training (resp. validation) subsets were merged. The performances on individual DMS datasets are then reported as the average of the 5 models trained on independent training sets.

Model design choices We experimented with mixture of diagonal Gaussian (MOG) priors with K = 1, 10, 100 mixture components. For K = 1, the mean and standard deviation of the prior distribution are fixed to $\mu_1 = \mathbf{0} \in \mathbb{R}^D$ and $\Sigma_1 = \text{diag}(0.01) \in \mathbb{R}^{D \times D}$. For K = 5, 10, 100, the trainable means are initialized randomly from a Normal distribution and for $k = 1, \dots, K$, $\Sigma_k = \text{diag}(0.01) \in \mathbb{R}^{D \times D}$ is fixed. We also experimented with a VAMP prior and K = 5, 10, 100mixture components. The prototypes in input space were initialized randomly and all trainable. In addition to experimenting with different prior distribution hyperparameters, we performed an ablation study of our model trained on DMS data. We report the performances of models with and without Transformer(.) and DwFC(.) layers. The different models are summarized in Table 1. Source Model Name Short Description matVAE-MSA Our Matrix variational auto-encoder trained on MSA data (Fig. 1) VAMPk matVAE-MSA with a Vamp prior and k components experiments MOGk matVAE-MSA with a MOG prior and k components Vec-DMS Encoder trained on DMS data, without Transformer or DwFC DwFC-DMS Encoder with DwFC trained on DMS data, without Transformer matENC-DMS Matrix encoder trained on DMS data (Fig. 2) Best performing model flavor on a given DMS dataset ProteinGym Best Benchmark VAE-based model DeepSequence ESM PLM-based model Tranception PLM and position-wise AA frequency-based model TranceptEVE PLM, EVE and position-wise AA frequency-based model MSA Transformer Position-wise transformer and PLM-based model GEMME Evolutionary tree based model VESPA Linear Ensemble of PLM, bio-mechanic features and positionwise frequency models. MIF PLM and inverse folding model Site-Independent Position-wise entropy-based model

Table 1: Summary of baselines and experiments.

Baselines All models with zero-shot prediction performances reported in the ProteinGym bench-342 mark were considered for inclusion as a baseline. Only those models that demonstrated the highest 343 SpearmanR zero-shot performances on at least one pharmacogene-related protein DMS dataset were 344 used as a baseline. These models fall into one of the following model families: DeepSequence (Ries-345 selman et al., 2018), ESM (Rives et al., 2021), Tranception (Notin et al., 2022a), MIF (Laine et al., 346 2019), GEMME (Laine et al., 2019), VESPA (Marquet et al., 2022), MSA Transformer (Rao et al., 347 2021). Within each family, the top-performing models vary in configuration (e.g., different param-348 eter counts), which we refer to as model "flavors" (See Table A.6). For example ESM2 (150M), 349 ESM2 (15B) and ESM-1v (ensemble) are all distinct flavors within the ESM family. We report the performances at the level of model families, by the average performances of the best performing 350 model flavors of that family on individual DMS datasets. The details of which model flavor performs 351 best on which DMS dataset are shown in Table A.5. We also compare with the "Site-Independent" 352 model of ProteinGym, which is similar to SIFT and Polyphen-2, both still widely used in pharma-353 cogenomics. In addition, we add a difficult baseline referred to as "Best Benchmark", which is the 354 best performing model flavor across all model families for each DMS dataset (See Table A.6). 355

To compare our models trained only on DMS data, we use the supervised learning baselines from 356 ProteinGym with 5 "Random" cross-validation splits. For this case the models belong to the follow-357 ing families: ESM, TranceptEVE, Tranception, DeepSequence, MSA Transformer. To the best of 358 the authors knowledge, all these baselines consist of zero-shot prediction models trained on MSA 359 data, and used pretrained in a supervised learning framework with embeddings of the protein se-360 quence (See (Notin et al., 2023a)). ProteinNPT is a slightly different architecture which jointly 361 trains on MSA and DMS data (Notin et al., 2023b). The details of the best performing model flavors 362 are in Table A.7. 363

364 The prediction baseline models are summarized in Table 1.

365 366

367 368

369

324

325

326

327

328

330

331

332

333

334

335

336

337 338

339 340 341

4 RESULTS & DISCUSSION

4.1 CHOICE OF THE PRIOR DISTRIBUTION

370 Our experiments on zero-shot prediction tasks show that all our models with MOG and Vamp prior 371 distributions have similar average SpearmanR and AUROC performances, regardless of the num-372 ber of components (Fig. 3a & Table 2c). All the models we evaluated perform slightly better than 373 MIF, Site-Independent and VESPA, but worse than all other baselines we chose (Table 2c). The 374 standard deviation across proteins of our model is rather large, similarly to the rest of the baseline 375 models, which prevents us from drawing strong conclusions. Numerically, our two best performing algorithms have a MOG prior with 1 and 10 components respectively, and similar reported average 376 SpearmanR of 0.401 and 0.400 respectively. At the dataset level, 2/26 proteins: TPOR and SCN5A, 377 have an increase of at least 10% in SpearmanR performances for at least one latent prior with more 378 than one component (Table 3b). A 10% threshold for relative improvement compared to MOG1 379 effectively separates outliers (Fig. A.9). Overall the choice of priors did not bring the expected im-380 provement in neither SpearmanR or AUROC performances. Numerically, the worst performing al-381 gorithm has a MOG prior with 5 components and an average SpearmanR of 0.395. For the rest of the 382 analysis, we use the model with a MOG prior and 1 component and refer to it as "matVAE-MSA". We choose this model because it is the model with the lowest complexity and which performs best on 383 average among our models trained on MSA data. Notably, for zero-prediction tasks, matVAE-MSA 384 outperforms the best benchmark model for two drug target proteins: RAF1 and MK01 (Table 2b). 385 These are two proteins involved in multiple cellular pathways, with several variants associated with 386 the Noonan syndrome (Motta et al., 2020). Improved predictions of variant effect for these proteins 387 can improve diagnosis and management of this syndrome, by reducing the effects of congenital heart 388 defect and several deformities (Pandit et al., 2007). 389



	MOGk			VAMPk		
UnitProt ID	5	10	100	5	10	100
TPOR (%)	10.3	32.6	16.8	22.7	24.4	1.2
SCN5A (%)	12.5	28.5	-0.8	-19.3	-23.4	32.5

(a) Raw performances. One point denotes a DMS dataset and the boxplots describe the distribution of scores across DMS datasets.

(b) Relative increase compared to MOG1. We show the 2 proteins for which the relative increase is greater than +10% for at least one choice of prior/number of components.

Figure 3: SpearmanR zero-shot performances for various choice of latent priors.

4.2 TRAINING ON DMS DATASETS

397

398

399

400 401

402 403

404

405 The protein specific models trained on DMS data (matENC-DMS) perform more than 25% better 406 than the similar model trained on MSA (matVAE-MSA) with respect to the average SpearmanR (Ta-407 ble A.4). Overall, all models expect Vec-DMS perform better than their zero-shot prediction counter 408 parts (Table 2a and Table 2c). The relative increase in performances of matENC-DMS compared to 409 matVAE-MSA across protein categories is shown in Fig. 4c. All the ADME related DMS datasets 410 show an increase in performances with matENC-DMS. Among "Drug target" proteins, 8 (10 DMS 411 datasets) show an increase in performances of more than +50% compared to matVAE-MSA. Since 412 the two models have similar capacity, one explanation could be the invalidity of the conservation assumption for these specific proteins. Among other potential confounders known by the authors 413 (e.g. quality and size of MSA or DMS, protein characteristics), none could individually explain the 414 differences in performances according to a correlation analysis (Fig. A.10). 415

416 Our supervised matENC-DMS model outperforms, for both average SpearmanR and AUROC 417 scores, all our chosen zero-shot prediction baselines, except "Best Benchmark", which include models with a lot more trainable parameters (Table A.4). Our ablation study shows that our model with-418 out both Transformer and DwFC layer (Vec-DMS) performs the worse overall (See also Fig. 4d). 419 Our model with DwFC layer and without Transformer performs much better, and slightly worse 420 than both ESM and matENC-DMS. This indicates that the transformer layer is not determinant 421 for the good results of matENC-DMS. This could be due to the design of transformer itself, or 422 the quality of the PDB structures predicted by AlphaFold which might put an inadequate induc-423 tive bias on the attention mechanism. Also, the models trained on DMS datasets have a larger 424 standard deviation than the rest of the baselines for both metrics. Fig. 4b shows graphically that 425 matENC-DMS outperforms the current zero-shot prediction "Best Benchmark" model on about half 426 (15/33) of the DMS datasets. In terms of the protein categories, the median SpearmanR perfor-427 mances of matENC-DMS are below "Best Benchmark" for both "Drug target" and "ADME-other", 428 and above for "ADME-transporter" and "ADME-CYP" (Fig. 4a). Further, the relative performances 429 of matENC-DMS versus "Best Benchmark" varies roughly between +75% and -75% for "Drug target" (Fig. 4b). This is likely explained by the diversity of proteins included in the "Drug Target" 430 category. The performances vary less than expected, roughly between +25% and -25%, for the 431 ADME related categories.



(c) Relative increase of matENC-DMS compared to matVAE-MSA.

(d) Relative increase of matENC-DMS compared to Vec-MSA.

Figure 4: SpearmanR performances of matENC-DMS versus other zero-shot prediction models, per protein category. One point denotes the score obtained on a DMS dataset and the boxplots describe 452 the distribution of scores across DMS datasets.

We show comparison with supervised learning baseline models in Table 2a. Vec-DMS performs the worse against all supervised learning baselines, while both matENC-DMS and DwFC-DMS perform better than DeepSequence on average, although the standard deviation is larger than other baselines. Overall our model perform similarly to the One-hot encoding method, but arguably worse than fine tuning approaches: TranceptEVE, MSA transformer, ESM, Tranception; and a joint training approach: ProteinNPT.

Table 2: Numerical results for our models against baselines. The models trained by us are in bold font.

Model name	SpearmanR
Best Benchmark	0.689 ± 0.162
ProteinNPT	0.671 ± 0.195
Tranception	0.646 ± 0.173
ESM	0.586 ± 0.153
MSATransformer	0.579 ± 0.17
TranceptEVE	0.536 ± 0.156
One-Hot Encoding	0.528 ± 0.168
matENC-DMS	0.522 ± 0.24
DwFC-DMS	0.507 ± 0.242
DeepSequence	0.501 ± 0.148
Vec-DMS	0.356 ± 0.236

(a) Numerical SpearmanR and AUROC ($\mu \pm \sigma$) for matENC-DMS with different architectures against baselines on supervised prediction tasks.

UniProt ID	matVAE-MSA	Best Benchmark
MK01	0.256	0.241
RAF1	0.541	0.482

shot prediction tasks.

Model Name	SpearmanR	AUROC
Best Benchmark	0.529 ± 0.151	0.794 ± 0.076
ESM	0.508 ± 0.157	0.78 ± 0.079
TranceptEVE	0.485 ± 0.167	0.763 ± 0.088
Tranception	0.478 ± 0.165	0.761 ± 0.086
GEMME	0.461 ± 0.155	0.750 ± 0.085
MSA Transformer	0.452 ± 0.167	0.746 ± 0.089
DeepSequence	0.425 ± 0.151	0.729 ± 0.084
MOG1	0.401 ± 0.138	0.721 ± 0.083
MOG10	0.400 ± 0.136	0.721 ± 0.084
VAMP100	0.399 ± 0.137	0.720 ± 0.082
VAMP10	0.396 ± 0.146	0.718 ± 0.087
MOG100	0.395 ± 0.142	0.717 ± 0.085
VAMP5	0.395 ± 0.14	0.716 ± 0.083
MOG5	0.395 ± 0.137	0.718 ± 0.084
MIF	0.394 ± 0.185	0.718 ± 0.091
Site-Independent	0.390 ± 0.145	0.713 ± 0.078
VESPA	0.383 ± 0.147	0.712 ± 0.091

(c) Numerical SpearmanR and AUROC ($\mu \pm \sigma$) for matVAE-MSA with different priors against baselines (b) SpearmanR for two proteins where matVAE-MSA on zero-shot prediction tasks. The table is sorted veroutperforms the Best Benchmark baseline for zerotically with respect to the SpearmanR score.

484 485

476

477

478 479 480

481

482

483

449

450

451

457

458

459

460

461

462 463

464

486 4.3 FUTURE WORK

488 Our results experimenting with expressive multi-modal priors did not show improvements compared
 489 to simple Gaussian prior. The investigations of the potential pitfalls of our current approach as well
 490 as further biology-relevant interpretations of the latent prior modes are left for future works.

491 Next, the joint use of both DMS and MSA data for training is an important next step towards im-492 proved model performances. In our work, our primary objective was to evaluate the information 493 provided by MSA and DMS data separately for variant effect prediction. DMS and MSA data could 494 nonetheless be used jointly for training, for instance in a fine tuning approach, where a VAE model 495 is first trained on MSA data, and the encoding part is fine tuned on DMS data. This might lead to 496 improved performances since both datasets would then contribute to the model performances. To-497 gether with experimentation on a wider range of proteins, fine tuning on DMS data is an interesting research direction that we leave as future work. 498

499 Lastly, the design of a model able to learn from multiple proteins is also an interesting next avenue 500 for research. Our current architecture can be extended to work for proteins of different lengths 501 L_1, L_2, \ldots This could be done by slightly modifying the DwFC layer (Eq. (4)). An idea would be to first define $L_M = \max(L_1, L_2, ...)$ and initialize $\mathbf{U} \in \mathbb{R}^{H \times L_M}$. At run time, with an input protein encoding $\underline{\mathbf{x}}' \in \mathbb{R}^{L \times d}$, the first L columns from matrix $\mathbf{U} \in \mathbb{R}^{H \times L_M}$ can be used, which 502 leads to: $\underline{s} = U_{(:,1:L)}\underline{x}' + \mathbf{b}$, where $U_{(:,1:L)}$ denotes the sub-matrix of U which includes all the 504 rows and the first L columns of U. The complexity of the function described in (4) depends on 505 the length of the protein under consideration, while the memory complexity scales linearly with the 506 length of the longest protein. This is the same as our current approach where one model is fit to 507 individual proteins. It differs in that the weights included up to a column l would be shared between 508 all the proteins of length at least l. This makes the transformer learn to organize information, by 509 placing what is relevant to all proteins in the first rows of the representation $\mathbf{x}' \in \mathbb{R}^{H \times d}$. An issue 510 with this approach when training on DMS data is that the meaning, support and distributions of the 511 DMS scores vary largely across DMS datasets, thus requiring the fitness scores to be standardized 512 (Fig. A.11). An interesting future research direction is the quantification of similarities in selection 513 assays of DMS datasets, so that they can exploited in regression models.

514 515

5 CONCLUSION

516 517

518 We proposed a transformer-based matrix variational auto-encoder and evaluated its performances 519 on DMS datasets of drug target and ADME proteins. We showed that advanced priors such as 520 mixture of Gaussians and VampPrior did not provide improvement over a simple Gaussian prior for the latent space when our model was trained on MSA data (matVAE-MSA). matVAE-MSA 521 nonetheless outperformed on average a model similar to widely used VEPs in pharmacogenomics. 522 Moreover, matVAE-MSA outperformed the best benchmark models in ProteinGym for 2 proteins 523 related to the Noonan syndrome. Our architecture allowed to compare performances with models 524 of similar capacity but trained on DMS datasets (matENC-DMS) instead of MSA. Although DMS 525 datasets were often much smaller, matENC-DMS outperformed matVAE-MSA for 15 out of 33 526 DMS datasets, including those of all ADME and certain drug target proteins. MSAs may thus be 527 limiting the performances of VEPs for some proteins for which the conservation assumption does not 528 hold. This in turn motivates the development of DMS datasets and the study of their relationships, 529 in order to further improve variant effect prediction.

530 531

532 REFERENCES

- Ivan A. Adzhubei, Steffen Schmidt, Leonid Peshkin, Vasily E. Ramensky, Anna Gerasimova, Peer
 Bork, Alexey S. Kondrashov, and Shamil R. Sunyaev. A method and server for predicting damaging missense mutations. *Nature Methods*, 7(4):248–249, April 2010. ISSN 1548-7105. doi: 10.1038/nmeth0410-248.
- 538

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer Normalization. *arXiv:1607.06450* [*cs, stat*], July 2016.

549

550

551

552

553

556

565 566

567

568

569 570

574

575

576 577

578

- 540
 541
 542
 542
 543
 543
 544
 545
 545
 546
 547
 548
 548
 549
 549
 549
 540
 541
 541
 541
 542
 543
 544
 544
 544
 544
 545
 546
 547
 548
 548
 549
 549
 549
 541
 541
 541
 541
 542
 543
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
- Nat Dilokthanakul, Pedro A. M. Mediano, Marta Garnelo, Matthew C. H. Lee, Hugh Salimbeni,
 Kai Arulkumaran, and Murray Shanahan. Deep Unsupervised Clustering with Gaussian Mixture
 Variational Autoencoders, January 2017.
- 548 R Durbin. Biological sequence analysis: Probabilistic models of proteins and nucleic acids, 1998.
 - J.-L. Durrieu, J.-Ph. Thiran, and F. Kelly. Lower and upper bounds for approximation of the Kullback-Leibler divergence between Gaussian Mixture Models. In 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4833–4836, March 2012. doi: 10.1109/ICASSP.2012.6289001.
- Douglas M. Fowler and Stanley Fields. Deep mutational scanning: A new style of protein science.
 Nature Methods, 11(8):801–807, August 2014. ISSN 1548-7105. doi: 10.1038/nmeth.3027.
- Jonathan Frazer, Pascal Notin, Mafalda Dias, Aidan Gomez, Joseph K. Min, Kelly Brock, Yarin
 Gal, and Debora S. Marks. Disease variant prediction with deep generative models of evolutionary data. *Nature*, 599(7883):91–95, November 2021. ISSN 1476-4687. doi: 10.1038/ s41586-021-04043-8.
- Sam Gelman, Sarah A. Fahlberg, Pete Heinzelman, Philip A. Romero, and Anthony Gitter. Neural networks to learn protein sequence–function relationships from deep mutational scanning data. *Proceedings of the National Academy of Sciences*, 118(48):e2104878118, November 2021. doi: 10.1073/pnas.2104878118.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity Mappings in Deep Residual Networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision* - ECCV 2016, pp. 630–645, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46493-0. doi: 10.1007/978-3-319-46493-0_38.
- 571 Chloe Hsu, Hunter Nisonoff, Clara Fannjiang, and Jennifer Listgarten. Learning protein fitness models from evolutionary and assay-labeled data. *Nature Biotechnology*, 40(7):1114–1122, July 2022. ISSN 1546-1696. doi: 10.1038/s41587-021-01146-5.
 - Po-Ssu Huang, Scott E. Boyken, and David Baker. The coming of age of de novo protein design. *Nature*, 537(7620):320–327, September 2016. ISSN 1476-4687. doi: 10.1038/nature19946.
 - M. Ingelman-Sundberg, S. Mkrtchian, Y. Zhou, and V.M. Lauschke. Integrating rare genetic variants into pharmacogenetic drug response predictions. *Human Genomics*, 12(1), 2018. ISSN 1473-9542. doi: 10.1186/s40246-018-0157-3.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03819-2.
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. arXiv:1312.6114 [cs, stat], May 2014.
- Elodie Laine, Yasaman Karami, and Alessandra Carbone. GEMME: A Simple and Fast Global
 Epistatic Model Predicting Mutational Effects. *Molecular Biology and Evolution*, 36(11):2604–2619, November 2019. ISSN 0737-4038. doi: 10.1093/molbev/msz179.

- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, March 2023. doi: 10.1126/science.ade2574.
- Céline Marquet, Michael Heinzinger, Tobias Olenyi, Christian Dallago, Kyra Erckert, Michael Bernhofer, Dmitrii Nechaev, and Burkhard Rost. Embeddings from protein language models predict conservation and variant effects. *Human Genetics*, 141(10):1629–1647, October 2022. ISSN 1432-1203. doi: 10.1007/s00439-021-02411-y.
- 604 Marialetizia Motta, Luca Pannone, Francesca Pantaleoni, Gianfranco Bocchinfuso, 605 Francesca Clementina Radio, Serena Cecchetti, Andrea Ciolfi, Martina Di Rocco, Mariet W. 606 Elting, Eva H. Brilstra, Stefania Boni, Laura Mazzanti, Federica Tamburrino, Larry Walsh, Kate-607 lyn Payne, Alberto Fernández-Jaén, Mythily Ganapathi, Wendy K. Chung, Dorothy K. Grange, Ashita Dave-Wala, Shalini C. Reshmi, Dennis W. Bartholomew, Danielle Mouhlas, Giovanna 608 Carpentieri, Alessandro Bruselles, Simone Pizzi, Emanuele Bellacchio, Francesca Piceci-609 Sparascio, Christina Lißewski, Julia Brinkmann, Ronald R. Waclaw, Quinten Waisfisz, Koen van 610 Gassen, Ingrid M. Wentzensen, Michelle M. Morrow, Sara Álvarez, Mónica Martínez-García, 611 Alessandro De Luca, Luigi Memo, Giuseppe Zampino, Cesare Rossi, Marco Seri, Bruce D. 612 Gelb, Martin Zenker, Bruno Dallapiccola, Lorenzo Stella, Carlos E. Prada, Simone Martinelli, 613 Elisabetta Flex, and Marco Tartaglia. Enhanced MAPK1 Function Causes a Neurodevelopmental 614 Disorder within the RASopathy Clinical Spectrum. American Journal of Human Genetics, 107 615 (3):499-513, September 2020. ISSN 1537-6605. doi: 10.1016/j.ajhg.2020.06.018. 616
- Pauline C. Ng and Steven Henikoff. SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Research*, 31(13):3812–3814, July 2003. ISSN 0305-1048. doi: 10.1093/nar/gkg509.
- Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena Hurtado, Aidan N. Gomez, Debora Marks, and Yarin Gal. Tranception: Protein Fitness Prediction with Autoregressive Transformers and Inference-time Retrieval. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 16990–17017. PMLR, June 2022a.
- Pascal Notin, Lood Van Niekerk, Aaron W. Kollasch, Daniel Ritter, Yarin Gal, and Debora S. Marks.
 TranceptEVE: Combining Family-specific and Family-agnostic Models of Protein Sequences for
 Improved Fitness Prediction, December 2022b.
- Pascal Notin, Aaron W. Kollasch, Daniel Ritter, Lood Van Niekerk, Steffan Paul, Han Spinner, Nathan J. Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko Franceschi, Yarin Gal, and Debora Susan Marks. ProteinGym: Large-Scale Benchmarks for Protein Fitness Prediction and Design. In *Thirty-Seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, November 2023a.
 - Pascal Notin, Ruben Weitzman, Debora Marks, and Yarin Gal. Proteinnpt: Improving protein property prediction and design with non-parametric transformers. *Advances in Neural Information Processing Systems*, 36:33529–33563, 2023b.
- Bhaswati Pandit, Anna Sarkozy, Len A. Pennacchio, Claudio Carta, Kimihiko Oishi, Simone Martinelli, Edgar A. Pogna, Wendy Schackwitz, Anna Ustaszewska, Andrew Landstrom, J. Martijn Bos, Steve R. Ommen, Giorgia Esposito, Francesca Lepri, Christian Faul, Peter Mundel, Juan P. López Siguero, Romano Tenconi, Angelo Selicorni, Cesare Rossi, Laura Mazzanti, Isabella Torrente, Bruno Marino, Maria C. Digilio, Giuseppe Zampino, Michael J. Ackerman, Bruno Dallapiccola, Marco Tartaglia, and Bruce D. Gelb. Gain-of-function RAF1 mutations cause Noonan and LEOPARD syndromes with hypertrophic cardiomyopathy. *Nature Genetics*, 39(8):1007–1012, August 2007. ISSN 1061-4036. doi: 10.1038/ng2073.
- 645 646

635

636

637

Roshan M. Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. MSA Transformer. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 8844–8856. PMLR, July 2021.

648	Adam J. Riesselman, John B. Ingraham, and Debora S. Marks. Deep generative models of genetic
649	variation capture the effects of mutations. <i>Nature Methods</i> , 15(10):816–822, October 2018, ISSN
650	1548-7105. doi: 10.1038/s41592-018-0138-4.
651	

- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings* of the National Academy of Sciences, 118(15):e2016239118, April 2021. doi: 10.1073/pnas. 2016239118.
- Borja Rodríguez Gálvez. An Information-Theoretic Approach to Generalization Theory. PhD thesis,
 KTH Royal Institute of Technology, 2024.
- Jakub M. Tomczak and Max Welling. VAE with a VampPrior, February 2018.
- Yasemin Bozkurt Varolgüneş, Tristan Bereau, and Joseph F. Rudzinski. Interpretable embeddings
 from molecular simulations using Gaussian mixture variational autoencoders. *Machine Learning: Science and Technology*, 1(1):015012, April 2020. ISSN 2632-2153. doi: 10.1088/2632-2153/
 ab80b7.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
 Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need, December 2017.
- Kevin K Yang, Niccolò Zanichelli, and Hugh Yeh. Masked inverse folding with sequence transfer for
 protein representation learning. *Protein Engineering, Design and Selection*, 36:gzad015, January
 2023. ISSN 1741-0126. doi: 10.1093/protein/gzad015.
- Yitian Zhou, Roman Tremmel, Elke Schaeffeler, Matthias Schwab, and Volker M. Lauschke. Challenges and opportunities associated with rare-variant pharmacogenomics. *Trends in Pharmacological Sciences*, 43(10):852–865, October 2022. ISSN 01656147. doi: 10.1016/j.tips.2022.07. 002.

702 A APPENDIX

704 A.1 DATASETS

722

723

724

706We train separate models on 26 pharmacogene-related proteins for which the DMS and MSA707datasets are readily available from the publicly available ProteinGym repository (Notin et al., 2023a).708The pharmacogene-related proteins were divided into four functional categories: Drug Targets709(n = 21) and Absorption-Distribution-Metabolism-Excretion (ADME) related proteins (n = 5).710The ADME category is further divided into Cytochrome ("CYP"), "transporter" and "other" ADME711proteins. In total we compare performances on 33 DMS datasets, some proteins having several DMS712datasets obtained under different selection assays (Table A.5).

713 A.1.1 PREPROCESSING OF MSA SEQUENCES

We followed the preprocessing steps proposed in DeepSequence for the MSA data (Riesselman et al., 2018). Sequences are removed from MSAs if they include more than 50% gaps. Columns are removed from a MSA if they contain more than 30% gaps across the MSA. For consistency, the columns that are removed from the MSA are also removed from the DMS datasets of that protein.
When training on MSAs, each sequence is sampled with a probability proportional to the reciprocal of the number of sequences within a given Hamming distance from that sequence (Riesselman et al., 2018). A summary description of the datasets is provided in Table A.3.

Table A.3: Description of DMS and MSA datasets per protein category. The most extreme values are in bold. L: Preprocessed sequence length; MSA Num Seq (resp. DMS Num Seq): number of sequences in the MSA (resp. DMS) datasets. ADME trans.: ADME Transporter.

Category		L	MSA Num Seq	DMS Num Seq
ADME CYP (n=2)	$\mu \pm \sigma$	490 ± 0	260849 ± 0	6256 ± 161
	min/max	490 / 490	260849 / 260849	6142 / 6370
ADME other (n=2)	$\mu \pm \sigma$	204 ± 57	86361 ± 94716	3246 ± 569
	min/max	164 / 245	19387 / 153335	2844 / 3648
ADME tran. (n=3)	$\mu \pm \sigma$	579 ± 44	144978 ± 90553	10491 ± 951
	min/max	553 / 630	40416 / 197259	9803 / 11576
Drug target (n=26)	$\mu \pm \sigma$	498 ± 427	62331 ± 125352	3374 ± 3134
	min/max	31 / 1863	911 / 611225	63 / 12464

A.2 MODEL ARCHITECTURES HYPER-PARAMETERS

736 The encoding and decoding transformers of matVAE-MSA are designed with 3 transformer layers 737 described in (3). The embedding dimensions of all the layers are identical and equal to d. The use of 738 3 layers allows to use information from neighbors up to order 3 according to the graph defined by the 739 attention mask. The attention mask is a thresholded distance matrix derived from the WT protein 740 structures predicted by Alphafold2 (Jumper et al., 2021). This means that queries are allowed to 741 attend to keys in the attention dot product if the predicted distance in 3d between the corresponding 742 AAs is $\leq c$. We chose c = 7Å which had the best performances for 4 out of 5 graph neural 743 network-based models predicting variant effects in (Gelman et al., 2021, Table S3). The structures 744 are readily available in the ProteinGym repository as PDB files (Notin et al., 2023a). For DwFC, 745 we chose $H = \min(H_{min}, L)$ with $H_{min} = 200$. This means that the dimension is not reduced for proteins with small enough sequence length $L \leq H_{min}$. Following the discussion in Section 746 2.1.1, we experimented with $H_{min} \approx d$ on some proteins, but could not get good performances. 747 For FCB, we used a 2-layer ReLU network with 1000 and 300 neurons, and output latent space 748 dimension D = 50, this is similar to the design of DeepSequence (Riesselman et al., 2018). The 749 networks f_{μ} and f_{σ} are both 1-layer FC ReLU networks with D neurons and output dimension 750 D. For stable computation of the closed form KLD, f_{σ} practically outputs $\ln \sigma^2$ using an output 751 pointwise operation: $x \mapsto \ln(\ln(e^x + 1))$. Symmetrical design choices were used for the decoding 752 part of matVAE-MSA. 753

For matENC-DMS, the *D*-dimensional latent vector in the FC bottleneck is passed into a prediction
 head with a 2-layer fully connected ReLU network, with 50 and 25 neurons, and an output dimension
 of 1. The decoding parts of matVAE-MSA are not used.

756 A.3 MODEL TRAINING

For matVAE-MSA, the models are trained on protein specific MSAs to minimize the negative ELBO in (9). The expected reconstruction error is approximated with a 1-sample Monte Carlo method. The reconstruction error is the cross-entropy between the true $\underline{\mathbf{x}} \in \mathbb{R}^{L \times d}$ and the reconstructed $\underline{\hat{\mathbf{x}}}$. For matENC-DMS, no variational formulation is used. The loss function is the mean squared error (MSE) between the true label $y \in \mathbb{R}$ and the reconstructed label \hat{y} .

For all our included proteins, the loss function for matVAE-MSA is optimized using the ADAM optimizer, with a fixed learning rate $\lambda = 8e - 5$, a batch size of 256 and 300,000 training steps. For matENC-DMS, the loss function is optimized with a fixed learning rate $\lambda = 1e - 4$, a batch size of and 100,000 training steps. The learning rates were chosen similar to the optimal one reported for a graph neural network model in (Gelman et al., 2021, Table S3). The batch sizes were chosen to obtain the most efficient use of our hardware. In matVAE-MSA, the memory footprint is mainly due to the attention matrices in the encoder and decoder transformer. We double the batch size for matENC-DMS compared to matVAE-MSA since matENC-DMS only has an encoder transformer.

A.4 PROTEINGYM BEST PERFORMING MODELS FOR A ZERO-SHOT PREDICTION TASK, VERSUS OUR MODELS FOR BOTH ZERO-SHOT PREDICTION AND SUPERVISED LEARNING TASKS.

Model Name	SpearmanR	AUROC
Best Benchmark	0.529 ± 0.151	0.794 ± 0.076
matENC-DMS	0.522 ± 0.24	0.784 ± 0.124
ESM	0.508 ± 0.157	0.780 ± 0.079
DwFC-DMS	0.507 ± 0.242	0.772 ± 0.125
TranceptEVE	0.485 ± 0.167	0.763 ± 0.088
Tranception	0.478 ± 0.165	0.761 ± 0.086
GEMME	0.461 ± 0.155	0.750 ± 0.085
MSA Transformer	0.452 ± 0.167	0.746 ± 0.089
DeepSequence	0.425 ± 0.151	0.729 ± 0.084
matVAE-MSA	0.401 ± 0.138	0.721 ± 0.083
MIF	0.394 ± 0.185	0.718 ± 0.091
Site-Independent	0.390 ± 0.145	0.713 ± 0.078
VESPA	0.383 ± 0.147	0.712 ± 0.091
Vec-DMS	0.356 ± 0.236	0.695 ± 0.123

Table A.4: SpearmanR and AUROC performances ($\mu \pm \sigma$). The table is sorted vertically with respect to SpearmanR. The name of the models trained by us is in bold font. The baseline



RAW SPEARMANR AND AUROC PERFORMANCES FOR ZERO-SHOT AND SUPERVISED A.5 LEARNING TASKS

Figure A.5: Raw Spearman Rank correlation coefficient performances on all DMS datasets for zero-shot prediction task. We display the performances of our models against the best performing models in ProteinGym for pharmacogene-related proteins. The datasets are sorted in decreasing "best benchmark" performances.



Figure A.6: AUROC performances on all DMS datasets for zero-shot prediction task. We display the performances of our models against the best performing models in ProteinGym for pharmacogenerelated proteins. The datasets are sorted in decreasing "best benchmark" performances.



Figure A.7: SpearmanR performances on all DMS datasets for a supervised learning task. We display the performances of our models against the best performing models in ProteinGym for pharmacogene-related proteins. The datasets are sorted in decreasing "best benchmark" performances.

A.6 PERFORMANCES WITH ASSAY SELECTION TYPE CATEGORY



Figure A.8: SpearmanR performances of matENC-DMS versus other models per Selection Type sub-groups.





Figure A.10: Univariate correlation analysis for potential confounders of the relative increase of matENC-DMS from matVAE-MSA. None of the considered confounders are significant to explain the relative increase from MSA. A multivariate correlation analysis was performed and did not show any significance (not shown). MSA Num Seqs: Number of variants in MSA; Probability of Loss of Function Intolerance: Genes with a pLI close to 1 are often associated with haploinsufficiency and dominant genetic diseases; Expected (resp. Observed) SNVs: Expected (resp. Observed) Single-nucleotide variant in each gene; Observed/Expected (o/e): Constrained genes have fewer observed variants than expected (low o/e) and are under a higher degree of selection than less constrained genes. DMS Num Seqs: Number of variants in DMS data;

A.8 ADDITIONAL DATASET INFORMATION



Figure A.11: Empirical kernel density estimates of the density function of DMS datasets target score. The support and range of each density function are standardized to [0,1] for visualization. The true support of the density functions is annotated on the right hand side. The densities are displayed with an increment of 1 along the y-axis.

÷	n	0	7	
4	υ	0	1	

	~	~	~	
ъ		ĸ	ч	
	\sim	~	~	

1091
1092
1093
1094
1095
1096
1097
1098
1099

1087	Category	UniProt	L	MSA	DMS	Selection	Best Benchmark - Flavor
1088		ID		Num Sea	Num Sea	Туре	
1089	ADME	CP2C9	490	260849	6142	Binding	ESM2 (150M)
1090	CYP					6	
1091					6370	Expression	ESM2 (650M)
1002	ADME	NUD15	164	153335	2844	Expression	MSA Transformer (ensemble)
1092	other			10007			
1093		TPMT	245	19387	3648	Expression	ESM-1v (ensemble)
1094	ADME	522A1	553	197259	9803	Expression	ESM-IV (ensemble)
1095	transporter				10094	Activity	ESM-1v (ensemble)
1096		SC6A4	630	40416	11576	Activity	MSA Transformer (ensemble)
1007	Drug target	ACE2	805	10865	2223	Binding	MIF
1097		ADRB2	413	201108	7800	Activity	GEMME
1098		BRCA1	1863	974	1837	Organismal	VESPA
1099						Fitness	
1100		CCR5	352	611225	6137	Binding	Tranception L
1101		CD19	556	1171	3761	Binding	MIF
1101		KCNF1	120	2104	2315	Activity	ESNI-IFI TrancentEVE I
1102		KCNEI	129	2104	2313	Expression	Tranception M no retrieval
1103		KCNH2	31	13900	200	Activity	Tranception M
1104		MET	287	184827	5393	Activity	MSA Transformer (ensemble)
1105		MK01	360	123422	6809	Organismal Fitness	DeepSequence (ensemble)
1106 1107		MTHR	656	4724	12464	Organismal Fitness	ESM2 (150M)
1108		NPC1	1278	6234	63	Activity	Tranception S
1109					637	Activity	MSA Transformer (ensemble)
1110		OTC	354	134484	1570	Activity	ESM-IF1
1111		PAI1	402	51792	5345	Activity	MIF-ST
		PPARG	505	39639	9576	Activity	ESM2 (15B)
1112		RAF1	648	9609	297	Organismal	MSA Transformer (single)
1113		SCNEA	20	40050	224	Fitness	ECM (1 (-in-1-)
1114		SCNSA	32	49959	224	Fitness	ESM-IV (single)
1115		SRC	536	37311	3366	Organismal	ESM-1v (ensemble)
1116					0070	Fitness	
1117					3372	Activity	ESM-1v (ensemble)
1110		TDV 1	2/3	21338	3181	Organismal	ESMI-IV (ensemble) ESM2 (15B)
1110			243	21330	5101	Fitness	
1119		TPOR	635	911	562	Organismal	MSA Transformer (single)
1120						Fitness	
1121		VKOR1	163	14425	697	Activity	ESM-1v (ensemble)
1122					2695	Expression	Tranception L

Table A.5: Deep Mutation Scanning datasets details retrieved from ProteinGym. Uniprot ID: Uni-versal protein identifier; L: Preprocessed sequence length; MSA Num Seq (resp. DMS Num Seq): number of sequences in the MSA (resp. DMS) datasets. Selection Type: DMS assay selection type. Best Benchmark - Flavor: best performing model flavor for zero-shot prediction tasks.

