

GENVARFORMER: PREDICTING GENE EXPRESSION FROM LONG-RANGE MUTATIONS IN CANCER

Anonymous authors

Paper under double-blind review

ABSTRACT

Distinguishing the rare “driver” mutations that fuel cancer progression from the vast background of “passenger” mutations in the non-coding genome is a fundamental challenge in cancer biology. A primary mechanism that non-coding driver mutations contribute to cancer is by affecting gene expression, potentially from millions of nucleotides away. However, existing predictors of gene expression from mutations are unable to simultaneously handle interactions spanning millions of base pairs, the extreme sparsity of somatic mutations, and generalize to unseen genes. To overcome these limitations, we introduce GenVarFormer (GVF), a novel transformer-based architecture designed to learn mutation representations and their impact on gene expression. GVF efficiently predicts the effect of mutations up to 8 million base pairs away from a gene by only considering mutations and their local DNA context, while omitting the vast intermediate sequence. Using data from 864 breast cancer samples from The Cancer Genome Atlas, we demonstrate that GVF predicts gene expression with 26-fold higher correlation across samples than current models. In addition, GVF is the first model of its kind to generalize to unseen genes and samples simultaneously. Finally, we find that GVF patient embeddings are more informative than ground-truth gene expression for predicting overall patient survival in the most prevalent breast cancer subtype, luminal A. GVF embeddings and gene expression yielded concordance indices of $0.706^{\pm 0.136}$ and $0.573^{\pm 0.234}$, respectively. Our work establishes a new state-of-the-art for modeling the functional impact of non-coding mutations in cancer and provides a powerful new tool for identifying potential driver events and prognostic biomarkers.

1 INTRODUCTION

The completion of the Human Genome Project marked a turning point in biology, revealing that less than 2% of the human genome encodes proteins, leaving the remaining 98%—the non-coding genome—a mystery (Piovesan et al.). Over time, it has become clear that much of this non-coding DNA plays critical roles in regulating gene expression, influencing key biological processes that shape cellular identity and state. Non-coding mutations can alter these regulatory mechanisms by disrupting transcription factor (TF) binding sites (Carrasco Pro et al.), modifying chromatin accessibility (Sundaram et al.), or affecting genome organization (Katainen et al.).

In the context of cancer, which is driven by selective pressures distinct from heritable human evolution, non-coding mutations are particularly intriguing. The vast majority of these mutations are likely to be “passengers,” with little impact on tumor fitness. However, a small subset of “driver” mutations are thought to confer a selective advantage and drive carcinogenesis. The challenge lies in distinguishing these rare drivers from the background noise of passengers (Carter et al.).

The increasing availability of whole-genome sequencing data from thousands of cancer samples (Priestley et al.; Aaltonen et al.; Sosinsky et al.) has created new opportunities to investigate the landscape of non-coding mutations in cancer. Non-coding driver mutations must have a biological effect in order to improve tumor fitness. As the primary function of non-coding DNA is gene regulation, it is likely that non-coding driver mutations affect gene expression. As a result, *modeling the effects of non-coding mutations can provide strong evidence for whether they are a driver and link their function to specific genes, revealing how they contribute to tumor fitness.*

054 Three challenges stand out when building models of gene expression in cancer. First, mutations that
055 affect gene expression can be millions of base pairs (Mbp) away from the gene they affect (Tjalsma
056 et al.). Second, somatic mutations can be relatively infrequent, with an average of 6 mutations per
057 Mbp in non-pediatric cancers (Poulsgaard et al.). Finally, because the space of possible mutations
058 is so large, almost all observed mutations are unique to each tumor. Prior work has tried to address
059 these challenges by using lasso models on so-called mutation “hotspots” (Zhang et al.; Soltis et al.;
060 Pudjihartono et al.). Using hotspots can often increase the frequency of features by one to two orders
061 of magnitude by aggregating mutations that are close together or are in a locus with more mutations
062 than expected by chance. While this addresses the first issue by omitting all DNA context, it also
063 necessitates fitting a model for each gene of interest. In non-cancerous samples, the issue of low
064 mutation frequency has been addressed by sequence models trained on common and rare germline
065 variants simultaneously (Drusinsky et al.; Rastogi et al.; Spiro et al.). However, these attempts have
066 been limited to considering variants at most 25 kilobases (kbp) away from the start of a gene since
067 they require contiguous DNA sequences as input.

068 To address these challenges, we introduce GenVarFormer (GVF), a novel method for learning rep-
069 resentations of mutations that affect gene expression in cancer (§3). GVF predicts gene expression
070 from mutations with a correlation across samples over 26 times greater than current approaches
071 (§4.2). Furthermore, to the best of our knowledge, GVF is the first model to generalize to pre-
072 dicting cancer gene expression in unseen genes and samples using mutations alone. We then use
073 GVF to compute patient embeddings and find that, in the luminal A subtype of breast cancer, GVF
074 embeddings are more informative for patient prognosis than ground truth gene expression (§4.3).

076 2 RELATED WORKS

079 **Models of frequently mutated regions** In bioinformatics, a common way to identify mutations
080 that affect gene expression in cancer is to predict gene expression from mutation hotspots using lasso
081 regression (Zhang et al.; Soltis et al.; Pudjihartono et al.). This approach is attractive since it can
082 efficiently incorporate mutations spanning millions of nucleotides and using hotspots can increase
083 the frequency of the input features. However, it suffers from several major limitations. First, it
084 is sensitive to the hotspot calling algorithm used and discards potentially informative non-hotspot
085 mutations. Second, these models don’t learn generalizable features, which prevents transfer learning
086 and task adaptation. Finally, it requires fitting models per gene, preventing application to novel genes
087 that can be critical cancer drivers, for example gene fusions (Dashi & Varjosalo) and recently evolved
088 *de novo* genes (Xiao et al.). GenVarFormer (GVF) eliminates these issues: it has no dependency on
089 hotspot calling algorithms, uses all available somatic mutations, generates informative embeddings
090 for downstream tasks, and is a pan-gene model that generalizes to unseen genes.

092 **DNA sequence-to-function models** DNA sequence-to-function models for gene expression,
093 which typically apply convolutional neural networks (CNNs) or transformers to one-hot encoded
094 DNA, have recently scaled to contexts of up to 1 Mbp and achieved high accuracy across
095 genes (Avsec et al.). However, these models are not trained on paired human genetic variation
096 and gene expression, and evaluations consistently find they fail to predict expression differences
097 between individuals (Huang et al.; Sasse et al.). While several groups have trained or fine-tuned
098 models specifically on genetic variation to address this gap (Drusinsky et al.; Rastogi et al.; Spiro
099 et al.), the computational expense has limited them to input contexts of only 49 kbp. This narrow
100 context window severely restricts model performance for two key reasons. First, it is too small to
101 capture sparse functional variation. For example, with an average of only 6 somatic mutations per
102 Mbp in cancer (Poulsgaard et al.), a 49 kbp window is unlikely to contain a relevant signal. In the
103 dataset we used, we find that for 88% of genes, a 49 kbp window contains no mutations, making
104 accurate prediction impossible. Second, no model trained on genetic variation has demonstrated
105 generalization to unseen genes. This is a critical failure, as the biological processes governing gene
106 regulation—such as transcription factor binding to promoters and enhancers—are fundamentally
107 shared across the human genome. GVF uses a DNA context window of 16 Mbp—over 340 times
larger than prior work—to model the impact of sparse genetic variation and generalizes to unseen
genes.

Vector representations of somatic mutations Several efforts have been made to learn representations of somatic mutations for tasks such as cancer type classification and patient survival prediction (Kim et al.; Gupta et al.; Sanjaya et al.; Anaya et al.). Most focus exclusively on coding mutations to either identify them as driver mutations or predict tumor-level phenotypes such as cancer type. For instance, Mut2Vec employs a word2vec-inspired model to embed genes based on their co-occurrence patterns in patient mutation profiles, evaluating whether the representations are informative for finding driver mutations (Kim et al.). Anaya et al. is the most related to our work by using a similar mutation input structure, but demonstrates a multiple-instance learning framework to predict tumor type and microsatellite status, both of which are tumor-level phenotypes. Like other methods, it is also restricted to coding mutations and does not incorporate variant allele frequency (VAF), which can serve as a proxy for what fraction of cells in a tumor have acquired the mutation (Castro et al.). Overall, these methods are less likely to reveal insights about the biological function of mutations because they are trained to predict tumor-level phenotypes that are too high-level to reflect specific biological processes. To this end, several of these methods also bin the genomic position of mutations to the nearest megabase, increasing the representational similarity of mutations but further obfuscating their function. In contrast, GVF is the first model designed to learn representations of non-coding somatic mutations by training on the fundamental biological task of predicting gene expression, dramatically increasing the number of instances the model can learn from and enabling it to learn functionally relevant representations.

3 GENVARFORMER

We formulate predicting gene expression from mutations as a regression task where each instance is a particular gene and sample. The inputs are mutations from a 16 Mbp window centered on the gene transcription start site (TSS) along with gene-identifying information. A schematic of the full architecture is shown in Fig. 1. We define mutations with the following properties:

- **ALT**: the DNA sequence of the mutation, the alternative allele relative to the reference genome.
- **ILEN**: the indel length of the mutation. Negative for deletions, zero for substitutions, and positive for insertions. All mutations are left normalized (Tan et al.) and atomized such that all substitutions are single-base substitutions and the first nucleotide of an indel corresponds to the reference genome.
- **VAF**: variant allele fraction, the proportion of sequencing reads that have the mutation in a sample. This provides information about intratumoral heterogeneity and the prevalence of the mutation in a sample.
- **Flanking DNA**: 32 bp from both the 5' and 3' ends of the mutation.
- **POS**: the position of the mutation relative to the length of the chromosome it is on.

We then encode each set of input mutations as a sequence of vectors $\mathbb{M} = \{m_0, \dots, m_n\} \in \mathbb{R}^d$. More specifically, for each mutation the ALT is encoded via a single layer transformer followed by mean pooling. The ILEN and VAF are stacked into a 2-dimensional vector and projected into \mathbb{R}^d . The flanking DNA sequences are concatenated, pass through a shallow convolutional neural network (CNN), and mean pooled; in our experiments we use a randomly initialized ConvNova (Bo et al.) module for this. This yields three d -dimensional vectors that are summed together to get each of the mutation vectors \mathbb{M} . Finally, to help prevent overfitting, each mutation's position is made relative to the gene's start position and rounded to the nearest hundreds place. The amount of rounding balances between making mutations less unique and the informativeness of a mutation's position relative to a gene. The positions of the mutations are used in subsequent transformer layers for rotary positional embeddings (Su et al.).

Gene specific features are also included in the input since only 6 mutations occur per Mbp on average (Poulsen et al.). If GVF only used mutations and their flanking DNA as input, this could make it impossible to identify the gene being predicted when mutations are mostly distal from the gene. In order to ensure these gene specific features are generalizable to unseen genes, we use features derived from DNA sequence. In practice, this can be the gene promoter and/or coding DNA, and in our experiments we use embeddings from Borzoi (Linder et al.), a sequence-to-function model.

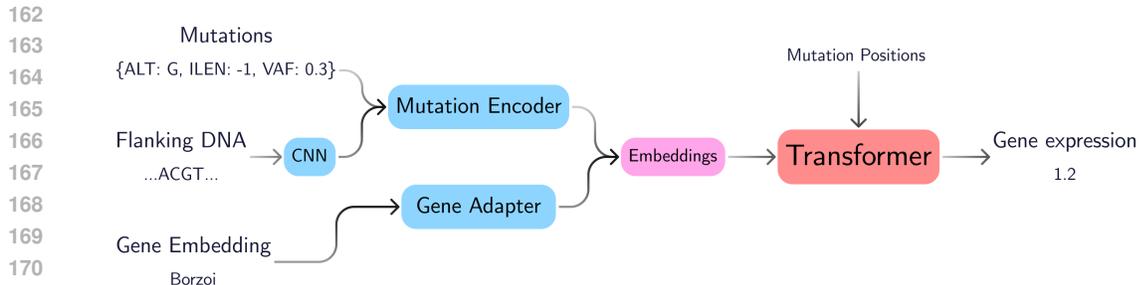


Figure 1: Schematic of GenVarFormer (GVF). Mutations from a 16 Mbp window around a gene are given as input. Each mutation’s information including the ALT, ILEN, VAF, and flanking DNA are transformed into an \mathbb{R}^d vector, as well as a gene embedding from Borzoi (Linder et al.) and all are summed together to get embeddings. These pass through a transformer along with mutation positions that are used for rotary positional embeddings (Su et al.). Finally, the embeddings for each mutation are mean pooled to a single vector and projected to predict gene expression.

As the gene adapter, we use a randomly initialized ConvNova (Bo et al.) module followed by mean pooling to compute a d -dimensional vector, which is added to all mutation vectors \mathbb{M} . In our case where gene embeddings are from a pretrained model, we multiply them by a learnable parameter initialized to a small value, e.g. 10^{-6} . This prevents the gene embedding from dominating the inputs at the start of training.

After encoding the mutations and gene-specific features into \mathbb{M} , they are passed through a transformer, mean-pooled, and projected to a scalar value to predict gene expression. For the loss function, we use gradient aligned regression (Zhu et al.) as minimizing pairwise distances has been shown to be critical in prior work predicting gene expression outside of cancer (Drusinsky et al.; Rastogi et al.; Spiro et al.). We additionally follow prior work and ensure that each batch seen during training exclusively consists of instances from the same gene.

3.1 TECHNICAL ADVANCES

Several technical challenges emerged while building GVF. First, we observed that the distribution of mutations per instance in the dataset was approximately Zipf-distributed (Fig. 2A). This meant that conventional padding strategies would lead to batches consisting of almost 100% pad tokens and make GVF infeasible for practical use (Fig. 2B). We thus implemented GVF using PyTorch (Paszke et al.) nested tensors to eliminate the need for padding during training and inference. However, even without padding we found that naive random sampling would cause out-memory-errors, as the number of mutations per batch was not limited by the batch size. This was remedied by implementing a bin packing sampler that ensured no batch contained more than a set number of mutations (Fig. 2C). As an added benefit, this sampling strategy also maximized the number of mutations per batch that could fit into memory and boosted GPU utilization. Finally, unlike natural language, mutations are not regularly spaced and to our knowledge, there is no implementation of rotary positional embeddings (RoPE) (Su et al.) that supports arbitrarily positioned tokens. We thus implemented a new RoPE Triton kernel for arbitrarily positioned tokens, building on the implementation from FlashAttention (Dao). Similarly, we used FlashAttention for all attention operations in GVF as it was the only implementation we found with a robust and performant forward and backward pass for nested tensors.

After overcoming these technical challenges, we were able to apply GVF to the full dataset without any issues. To demonstrate the gains in efficiency of using GVF over conventional biological sequence models, we benchmarked the throughput of GVF against Flashzoi (Hingerl et al.), a FlashAttention-enhanced version of Borzoi (Linder et al.), with matched window sizes of 524,288 bp. For GVF, we set the maximum number of variants per batch to 32,768, and for Flashzoi we used a batch size of 8—the largest power of 2 that would fit into GPU memory. We found that GVF was over 1,170 times faster than Flashzoi while using less GPU memory (Fig. 2D).

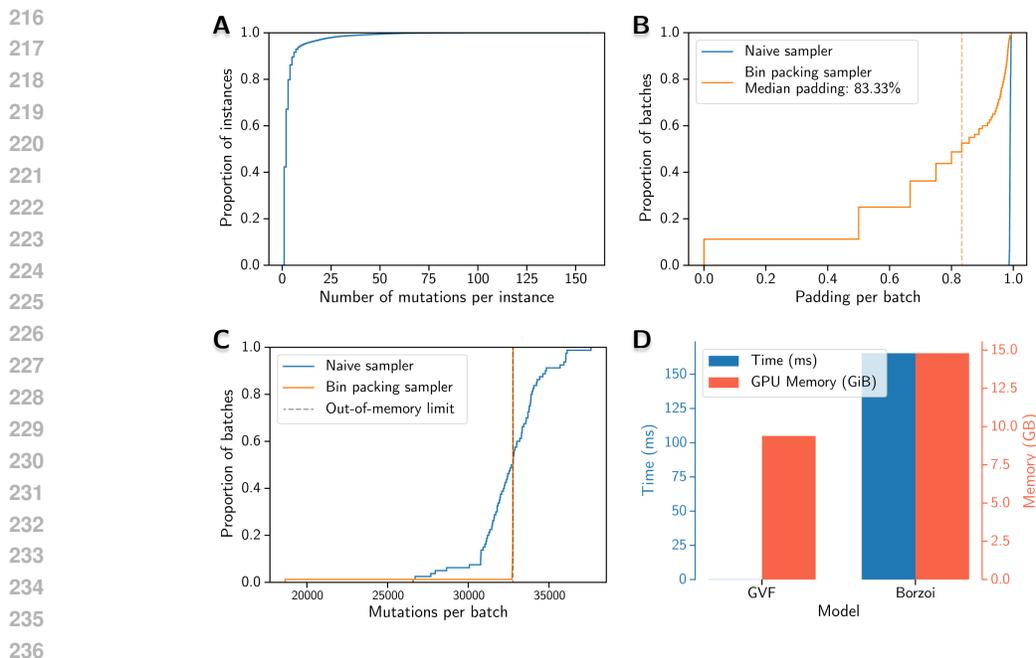


Figure 2: Motivation and benefit of using nested tensors throughout GVF. A) The distribution of mutations per instance in the dataset. This is approximately Zipf-distributed. B) The Zipf-distribution of mutations per instance causes extremely high amounts of padding per batch, even with the bin packing sampler. Dashed orange line is the median padding per batch with the bin packing sampler. C) Naive random sampling with a fixed batch size leads to out-of-memory errors since the number of mutations in a batch varies. The bin packing sampler strictly respects memory limits and increases GPU utilization by maximizing the the number of mutations per batch. D) Comparison of the average inference time per instance of GVF and Borzoi, as well as their peak memory usage. GVF is over 1,170 times faster at computing predictions.

4 EXPERIMENTS

We conduct two main experiments with GenVarFormer (GVF). First, we trained and evaluated it for predicting gene expression in tumors from somatic mutations, assessing its generalization to unseen samples, unseen genes, and simultaneously unseen samples and genes. Second, we computed patient embeddings with GVF and used linear probes to assess their clinical utility by predicting patient progression-free and overall survival, PAM50 subtype (Parker et al.), and early vs. late tumor stage.

4.1 DATA, SPLITTING, AND TRAINING

We obtained paired whole-genome and bulk RNA sequencing for 864 breast cancer samples from The Cancer Genome Atlas (TCGA) (Weinstein et al.). Since solid tumor biopsies are almost always a mixture of cancer and non-cancer cells, bulk RNA-seq measures a blend of cancer and non-cancer cells. This confounds the task of predicting cancer gene expression from mutations. Using the bulk expression encourages the model to learn how mutations in cancer cells affect gene expression in other cell types, which is mediated by dramatically different biological processes than within-cell gene regulation. To focus on cancer gene regulation, we used InstaPrism (Hu & Chikina) to estimate the cancer gene expression (*in silico* purification), removing the amount attributable to other cell types. We then generated splits of unseen samples (US), unseen genes (UG), and both (USG) for testing, validation, and training. We split the genes by chromosome such that the test and validation splits each had approximately 10% of the genes with any non-zero measurements. We then used 5-fold cross-validation across samples within the training split for hyperparameter tuning. Similar to prior work (Zhang et al.; Soltis et al.; Pudjihartono et al.), we finally regressed out the top 10 principal components of gene expression to remove batch effects and z-scored the residuals, respecting the 5 training folds, validation, and test splits. For training, we subset the dataset to genes

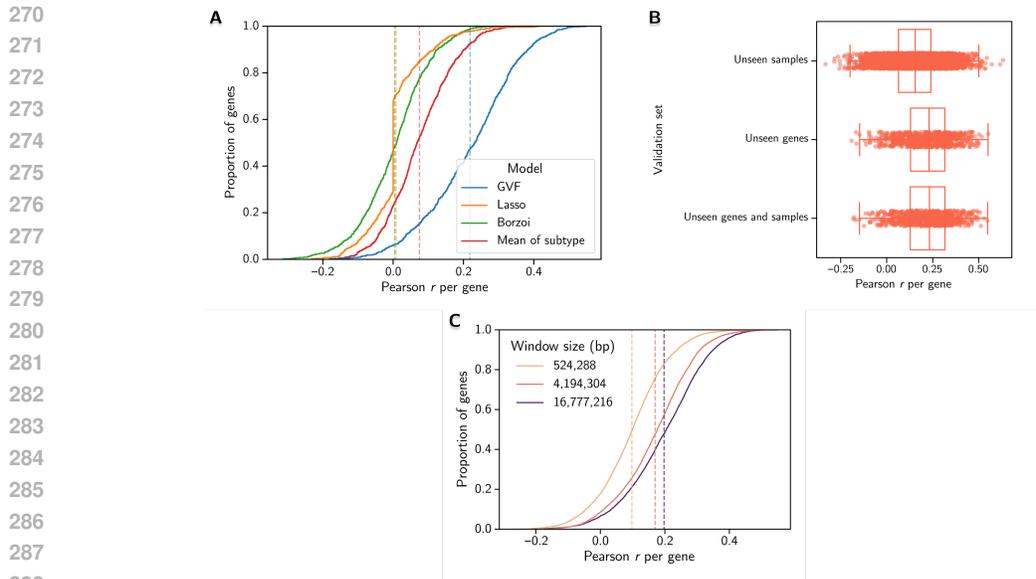


Figure 3: A) Performance of GVF, lasso models using mutation hotspots, Borzoi, and the mean gene expression of each breast cancer subtype. GVF performs over 26 times better than lasso models using the same input modality. Dashed, colored lines are the mean of each model’s performance. B) GVF demonstrates generalization in all three validation sets, most notably in unseen genes and samples with an average Pearson r of 0.219. C) As a result of improved scalability, GVF was trained with 16 Mbp windows which yielded double the performance of 524 kbp windows.

with a mean, purified expression greater than 1 (in $\log(\text{TPM} + 1)$). After training to predict gene expression, we used 3-fold nested cross validation for benchmarking on downstream tasks.

4.2 PREDICTING GENE EXPRESSION

To benchmark GVF, we computed the performance of three baselines: lasso models using hotspots called with the algorithm described in Zhang et al., predictions from Borzoi (Linder et al.) for all breast and breast cancer cell lines it was trained on, and the mean expression of the training data for each gene and subtype. Borzoi outputs 32-bp resolution tracks, so we follow (Linder et al.) and summed the predictions across exons to compute gene-level predictions. We did not fine-tune Borzoi due to computational constraints. We quantified performance by computing the Pearson correlation across samples for each gene and then computing the average of each gene’s correlation. We found that Borzoi yielded the lowest average correlation at 0.0043, followed by the lasso models at 0.0081, the mean subtype at 0.0749, and GVF at 0.2187. Borzoi’s low performance is unsurprising given that it was never trained on genetic variation and this result is consistent with reports evaluating similar models in non-cancerous tissue (Huang et al.; Sasse et al.). The lasso models represent a typical approach taken in literature (Zhang et al.; Soltis et al.; Pudjihartono et al.), which GVF outperforms by over 26-fold (Fig. 3A). Notably, GVF achieves this while generalizing to unseen genes and samples (Fig. 3B), which has yet to be reported for any existing biological sequence model (Drusinsky et al.; Rastogi et al.; Spiro et al.). We include the performance of the average expression per gene and subtype from the training data as a strong but simple baseline to evaluate whether GVF could be learning this straightforward relationship between expression and subtype. Finally, we conducted an ablation study with window sizes of 2^{19} , 2^{22} , and 2^{24} , finding that GVF improved in performance with increasing window size, achieving roughly double the performance with a 16 Mbp vs. 524 kbp window (Fig. 3C).

4.3 EVALUATING GENVARFORMER REPRESENTATIONS FOR CLINICAL UTILITY

To evaluate whether GVF could generate clinically informative patient-level representations, we first extracted embeddings from every layer in the model for every instance in the dataset. Each instance

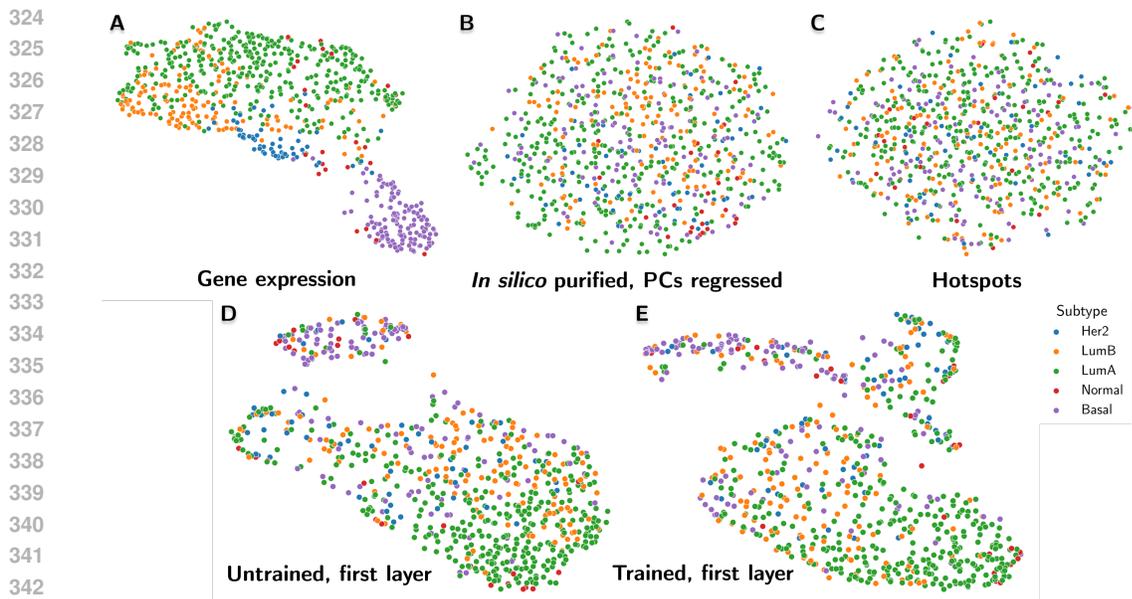


Figure 4: GVF patient representations display structure not seen in either outputs (B) or inputs (approximated by C). Tumors colored by subtype and projected via UMAP of A) gene expression, B) *in silico* purified gene expression with the top 10 PCs regressed, C) mutation hotspots, D) random projections from the mutation encoder/first layer of an untrained GVF model, and E) embeddings from the first layer of the trained GVF model. Random projections (D) are expected to show high-level patterns of mutations as a form of dimensionality reduction.

and layer yielded as many d -dimensional vectors as there were mutations, so we mean pooled the mutations to get gene-level embeddings. Finally, we concatenated these gene embeddings to get patient-level embeddings. We then generated UMAPs of gene expression, purified expression with the top 10 PCs regressed (the target that GVF was trained to predict), mutation hotspots, and embeddings from the mutation encoder/first layer of an untrained and trained GVF model (Fig. 4). As a case study, we then colored these projections by PAM50 (Parker et al.) subtype, an important prognostic biomarker for breast cancer. PAM50 subtypes are defined by a nearest-centroid classifier, Prediction Analysis of Microarray (PAM) (Tibshirani et al.), using the expression of 50 genes that were chosen to maximize progression-free risk stratification. Thus, it is expected that gene expression would largely reflect PAM50 subtypes as in Fig. 4A. After *in silico* purification and regressing out the first 10 principal components, subtype-specific patterns in gene expression are largely lost (Fig. 4B). Likewise, mutation hotspots do not appear to strongly correspond to subtype, nor show much overall structure (Fig. 4C).

Since breast cancer subtypes display distinct patterns of mutations (Perry et al.), we hypothesized that random projections of patients from an untrained GVF mutation encoder/first layer would also display some degree of clustering by subtype. The UMAP of patient embeddings from an untrained mutation encoder suggest this is true, as patients roughly stratify by subtype (Fig. 4D). This also shows that even random projections from GVF’s mutation encoder encode substantially more information than hotspots. Note that the random projections from the untrained mutation encoder only incorporate information about the DNA sequence resulting from mutations, not the positions, substitutions, or indels that occurred. The UMAP of GVF’s trained mutation encoder showed a similar, but finer-grained, substructure of patients (Fig. 4E). We quantified the informativeness of hotspots, random patient projections, and trained patient embeddings in the next experiment.

After this case study of PAM50 subtypes in the latent space, we focused on fitting linear probes to predict clinical annotations: progression-free and overall survival as well as PAM50 subtype and early vs. late stage by binning tumor stage into I-II and III-IV 1. We also fit and evaluated survival within each subtype since the two are strongly linked; subtype is a key biomarker that helps guide

378 treatment. We used mutation hotspots and an untrained GVF (i.e. no pre-training) as baselines
 379 for direct comparison to GVF. As a point of reference, we also evaluated advantaged feature sets:
 380 gene expression, *in silico* purified gene expression, and gene copy number. These features are not
 381 available to GVF and take advantage of either having RNA-seq available or knowledge of gene
 382 coordinates. As a result these feature sets also have 64 times fewer dimensions, corresponding to
 383 the dimensionality of GVF gene embeddings. For each combination of feature set and survival
 384 task, we projected the features onto their principal components and fit a Cox proportional hazard
 385 model to the projections. The number of principal components—and best layer of the model when
 386 applicable—was selected using the inner cross validation loop. For the non-survival tasks, we fit
 387 logistic regression models with an L2 penalty selected by cross validation. For several tasks, we
 388 observed that no mutation-based feature set yielded a score that was more than 1 standard deviation
 389 away from random, suggesting that with only 864 samples there may not have been enough data to
 390 fit meaningfully performant models. However, trained GVF embeddings were the most performant
 391 for every task where we could fit a mutation-based model that was substantially better than random.

392
 393 Table 1: Performance for predicting overall (OS) and progression-free survival (PFS), PAM50 sub-
 394 type, and early/late stage cancer. Early/late is defined as Stage I-II and Stage III-IV. Survival was
 395 also evaluated in each subtype. Values and standard deviations for survival are concordance indices,
 396 and for classification are for the area under the receiver-operating characteristic, which is one-vs-rest
 397 macro averaged for PAM50 subtype. Tasks where no mutation-only feature sets yielded performance
 398 greater than 1 standard deviation away from random performance were omitted. Advantaged fea-
 399 tures generally have a much higher signal-to-noise ratio (SNR) than mutations and are included for
 400 reference. Bold and underlined entries indicate the best and second-best non-random scores among
 401 mutation-only feature sets. Expr: gene expression. Pure: *in silico* purified expression. CN: gene
 402 copy number. Hotspots: mutation hotspots. GVF-R: randomly initialized, untrained GVF. PAM50:
 403 whether the features were subset to the PAM50 genes or not.

Features	OS:LumA	OS:Basal	PFS:LumA	PAM50	Early/Late
Sample size	$n = 437$	$n = 155$	$n = 437$	$n = 864$	$n = 864$
<i>Advantaged Features</i>					
Expr	0.53 \pm 0.28	0.77 \pm 0.18	0.63 \pm 0.1	0.98 \pm 0.0	0.61 \pm 0.04
Expr, PAM50	0.57 \pm 0.23	0.54 \pm 0.33	0.51 \pm 0.11	0.99 \pm 0.0	0.56 \pm 0.04
Pure	0.57 \pm 0.21	0.54 \pm 0.3	0.55 \pm 0.15	0.99 \pm 0.0	0.63 \pm 0.02
Pure, PAM50	0.42 \pm 0.28	0.64 \pm 0.3	0.55 \pm 0.13	0.98 \pm 0.0	0.54 \pm 0.07
CN	0.51 \pm 0.2	0.45 \pm 0.28	0.51 \pm 0.19	0.9 \pm 0.01	0.58 \pm 0.03
CN, PAM50	0.58 \pm 0.21	0.8 \pm 0.18	0.43 \pm 0.12	0.88 \pm 0.01	0.55 \pm 0.03
<i>Mutations Only (fair comparison to GVF)</i>					
Hotspots	0.47 \pm 0.24	0.31 \pm 0.26	0.46 \pm 0.09	0.61 \pm 0.03	0.52 \pm 0.03
GVF-R	0.64 \pm 0.22	0.43 \pm 0.29	0.66 \pm 0.16	0.72 \pm 0.02	0.5 \pm 0.05
GVF-R, PAM50	0.63 \pm 0.15	0.68 \pm 0.22	0.61 \pm 0.13	0.67 \pm 0.02	0.55 \pm 0.03
GVF	0.58 \pm 0.17	0.53 \pm 0.27	0.67 \pm 0.14	0.82 \pm 0.02	0.56 \pm 0.01
GVF, PAM50	0.71 \pm 0.14	0.71 \pm 0.21	0.65 \pm 0.14	<u>0.78 \pm 0.03</u>	<u>0.56 \pm 0.03</u>

424 5 CONCLUSION

425 We have presented GenVarFormer (GVF), a transformer-based model that expands our ability to
 426 functionally interpret non-coding mutations in cancer by predicting their effect on gene expression.
 427 By selectively modeling only mutations and their local sequence context, GVF efficiently models
 428 cancer gene regulation across up to 16 million base pairs. This context window is 340 times longer
 429 than the maximum length used by current personalized sequence-to-function models (Drusinsky
 430 et al.; Rastogi et al.; Spiro et al.) and runs over 1,170 times faster than a state-of-the-art sequence
 431

432 model, Flashzoi (Hingerl et al.). GVF also demonstrates a 26-fold increase in correlation across
 433 samples compared to traditional hotspot-based methods and is the first model of its kind capable
 434 of generalizing across both unseen genes and patients simultaneously. Finally, GVF patient em-
 435 beddings proved more informative for predicting survival in luminal A breast cancer than the gene
 436 expression data used for training.

437 One of the most important challenges in cancer is to identify the driver mutations that boost can-
 438 cer cells' ability to grow, evade the immune system, and otherwise improve their evolutionary
 439 fitness (Hanahan). As fundamental contributors to tumor fitness, driver mutations are invaluable
 440 biomarkers for patient prognosis and treatment and drug target discovery (Ostroverkhova et al.).
 441 Non-coding driver mutations must have an effect relevant to tumor fitness, and the primary function
 442 of non-coding DNA is gene regulation. Models of gene expression from mutations can therefore
 443 help discern which mutations meet this condition to be a driver mutation. Several challenges im-
 444 pede this task, as non-coding mutations can be multiple megabases away from the gene(s) they affect
 445 and they occur with very low frequency across sequence length and patients. GVF overcomes these
 446 challenges, offering a powerful new tool to quantify the consequences of non-coding mutations on
 447 gene expression and, as a result, prioritize non-coding driver mutations.

448 Several future directions stand out. Germline variants are known to influence cancer gene expres-
 449 sion (Li et al.), so incorporating them may improve predictive performance and reveal novel inter-
 450 actions between germline and somatic mutations. Extending to germline variants would also enable
 451 applications beyond cancer, where GVF may enable a more precise understanding of rare variants
 452 in genetic disease. We also did not delve into coding mutations in this work, and Anaya et al.'s find-
 453 ings suggest that indicating each mutation's position in any overlapping reading frames would be
 454 necessary to enable discrimination between classes of coding mutations. Despite these limitations,
 455 our model established a new state-of-the-art for predicting gene expression from mutations in can-
 456 cer. GenVarFormer helps to provide a path toward representing patient genomes more holistically,
 457 moving beyond a narrow focus on recurrent hotspots and coding mutation biomarkers.

458 6 ETHICS

459 Real patient data was used for this study. As such we conducted all work consistent with the data
 460 access policies set by the data distributor, Genomic Data Commons.

464 7 REPRODUCIBILITY

465 To ensure transparency and reproducibility, model code and weights will be publicly released.
 466 Datasets used for this work are available at the Genomic Data Commons Portal.

469 REFERENCES

470
 471 Lauri A. Aaltonen, Federico Abascal, Adam Abeshouse, Hiroyuki Aburatani, David J. Adams, Nis-
 472 hant Agrawal, Keun Soo Ahn, Sung-Min Ahn, Hiroshi Aikata, Rehan Akbani, Kadir C. Akdemir,
 473 Hikmat Al-Ahmadie, Sultan T. Al-Sedairy, Fatima Al-Shahrou, Malik Alawi, Monique Albert,
 474 Kenneth Aldape, Ludmil B. Alexandrov, Adrian Ally, Kathryn Alsop, Eva G. Alvarez, Fer-
 475 nanda Amary, Samirkumar B. Amin, Brice Aminou, Ole Ammerpohl, Matthew J. Anderson,
 476 Yeng Ang, Davide Antonello, Pavana Anur, Samuel Aparicio, Elizabeth L. Appelbaum, Yasuhito
 477 Arai, Axel Aretz, Koji Arihiro, Shun-ichi Ariizumi, Joshua Armenia, Laurent Arnould, Sylvia
 478 Asa, Yassen Assenov, Gurnit Atwal, Sietse Aukema, J. Todd Auman, Miriam R. R. Aure, Philip
 479 Awadalla, Marta Aymerich, Gary D. Bader, Adrian Baez-Ortega, Matthew H. Bailey, Peter J. Bai-
 480 ley, Miruna Balasundaram, Saianand Balu, Pratiti Bandopadhyay, Rosamonde E. Banks, Stefano
 481 Barbi, Andrew P. Barbour, Jonathan Barenboim, Jill Barnholtz-Sloan, Hugh Barr, Elisabet Bar-
 482 rera, John Bartlett, Javier Bartolome, Claudio Bassi, Oliver F. Bathe, Daniel Baumhoer, Prashant
 483 Bavi, Stephen B. Baylin, Wojciech Bazant, Duncan Beardsmore, Timothy A. Beck, Sam Behjati,
 484 Andreas Behren, Beifang Niu, Cindy Bell, Sercan Beltran, Christopher Benz, Andrew Berchuck,
 485 Anke K. Bergmann, Erik N. Bergstrom, Benjamin P. Berman, Daniel M. Berney, Stephan H.
 Bernhart, Rameen Beroukhim, Mario Berrios, Samantha Bersani, Johanna Bertl, Miguel Betan-
 court, Vinayak Bhandari, Shriram G. Bhosle, Andrew V. Biankin, Matthias Bieg, Darell Bigner,

- 486 Hans Binder, Ewan Birney, Michael Birrer, Nidhan K. Biswas, Bodil Bjerkehagen, Tom Boden-
 487 heimer, Lori Boice, Giada Bonizzato, Johann S. De Bono, Arnaud Boot, Moiz S. Bootwalla,
 488 Ake Borg, Arndt Borkhardt, Keith A. Boroevich, Ivan Borozan, Christoph Borst, Marcus Bosen-
 489 berg, Mattia Bosio, Jacqueline Boulton, Guillaume Bourque, Paul C. Boutros, G. Steven Bova,
 490 David T. Bowen, Reanne Bowlby, David D. L. Bowtell, Sandrine Boyault, Rich Boyce, Jeffrey
 491 Boyd, Alvis Brazma, Paul Brennan, Daniel S. Brewer, Arie B. Brinkman, Robert G. Bristow,
 492 Russell R. Broadus, Jane E. Brock, Malcolm Brock, Annegien Broeks, Angela N. Brooks,
 493 Denise Brooks, Benedikt Brors, Søren Brunak, Timothy J. C. Bruxner, Alicia L. Bruzos, Alex
 494 Buchanan, Ivo Buchhalter, Christiane Buchholz, Susan Bullman, Hazel Burke, Birgit Burkhardt,
 495 Kathleen H. Burns, John Busanovich, Carlos D. Bustamante, Adam P. Butler, Atul J. Butte,
 496 Niall J. Byrne, Anne-Lise Børresen-Dale, Samantha J. Caesar-Johnson, Andy Cafferkey, Declan
 497 Cahill, Claudia Calabrese, Carlos Caldas, Fabien Calvo, Niedzica Camacho, Peter J. Campbell,
 498 Elias Campo, Cinzia Cantù, Shaolong Cao, Thomas E. Carey, Joana Carlevaro-Fita, Rebecca
 499 Carlsen, Ivana Cataldo, Mario Cazzola, Jonathan Cebon, Robert Cerfolio, Dianne E. Chadwick,
 500 Dimple Chakravarty, Don Chalmers, Calvin Wing Yiu Chan, Kin Chan, Michelle Chan-Seng-
 501 Yue, Vishal S. Chandan, David K. Chang, Stephen J. Chanock, Lorraine A. Chantrill, Aurélien
 502 Chateigner, Nilanjan Chatterjee, Kazuaki Chayama, Hsiao-Wei Chen, Jieming Chen, Ken Chen,
 503 Yiwen Chen, Zhaohong Chen, Andrew D. Cherniack, Jeremy Chien, Yoke-Eng Chiew, Suet-
 504 Feung Chin, Juok Cho, Sunghoon Cho, Jung Kyoong Choi, Wan Choi, Christine Chomienne,
 505 Zechen Chong, Su Pin Choo, Angela Chou, Angelika N. Christ, Elizabeth L. Christie, Eric Chuah,
 506 Carrie Cibulskis, Kristian Cibulskis, Sara Cingarlini, Peter Clapham, Alexander Claviez, Sean
 507 Cleary, Nicole Cloonan, Marek Cmero, Colin C. Collins, Ashton A. Connor, Susanna L. Cooke,
 508 Colin S. Cooper, Leslie Cope, Vincenzo Corbo, Matthew G. Cordes, Stephen M. Corder, Isidro
 509 Cortés-Ciriano, Kyle Covington, Prue A. Cowin, Brian Craft, David Craft, Chad J. Creighton,
 510 Yupeng Cun, Erin Curley, Ioana Cutcutache, Karolina Czajka, Bogdan Czerniak, Rebecca A.
 511 Dagg, Ludmila Danilova, Maria Vittoria Davi, Natalie R. Davidson, Helen Davies, Ian J. Davis,
 512 Brandi N. Davis-Dusenbery, Kevin J. Dawson, Francisco M. De La Vega, Ricardo De Paoli-
 513 Iseppi, Timothy Defreitas, Angelo P. Dei Tos, Olivier Delaneau, John A. Demchok, Jonas De-
 514 meulemeester, German M. Demidov, Deniz Demircioğlu, Nening M. Dennis, Robert E. Den-
 515 roche, Stefan C. Dentre, Nikita Desai, Vikram Deshpande, Amit G. Deshwar, Christine Desmedt,
 516 Jordi Deu-Pons, Noreen Dhalla, Neesha C. Dhani, Priyanka Dhingra, Rajiv Dhir, Anthony DiBi-
 517 ase, Klev Diamanti, Li Ding, Shuai Ding, Huy Q. Dinh, Luc Dirix, Harsha Vardhan Doddapaneni,
 518 Nilgun Donmez, Michelle T. Dow, Ronny Drapkin, Oliver Drechsel, Ruben M. Drews, Serge
 519 Serge, Tim Dudderidge, Ana Dueso-Barroso, Andrew J. Dunford, Michael Dunn, Lewis Jonathan
 520 Dursi, Fraser R. Duthie, Ken Dutton-Regester, Jenna Eagles, Douglas F. Easton, Stuart Edmonds,
 521 Paul A. Edwards, Sandra E. Edwards, Rosalind A. Eeles, Anna Ehinger, Juergen Eils, Roland Eils,
 522 Adel El-Naggar, Matthew Eldridge, Kyle Ellrott, Serap Erkek, Georgia Escaramis, Shadrielle
 523 M. G. Espiritu, Xavier Estivill, Dariush Etemadmoghadam, Jorunn E. Eyfjord, Bishoy M. Fal-
 524 tas, Daiming Fan, Yu Fan, William C. Faquin, Claudiu Farcas, Matteo Fassan, Aquila Fatima,
 525 and The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analy-
 526 sis of whole genomes. 578(7793):82–93. ISSN 1476-4687. doi: 10.1038/s41586-020-1969-6.
 527 URL <https://www.nature.com/articles/s41586-020-1969-6>. Publisher: Na-
 528 ture Publishing Group.
- 526 Jordan Anaya, John-William Sidhom, Faisal Mahmood, and Alexander S. Baras. Multiple-instance
 527 learning of somatic mutations for the classification of tumour type and the prediction of mi-
 528 crosatellite status. 8(1):57–67. ISSN 2157-846X. doi: 10.1038/s41551-023-01120-3. URL
 529 <https://www.nature.com/articles/s41551-023-01120-3>.
- 530 Žiga Avsec, Natasha Latysheva, Jun Cheng, Guido Novati, Kyle R. Taylor, Tom Ward, Clare By-
 531 croft, Lauren Nicolaisen, Eirini Arvaniti, Joshua Pan, Raina Thomas, Vincent Dutordoir, Mat-
 532 teo Perino, Soham De, Alexander Karollus, Adam Gayoso, Toby Sargeant, Anne Mottram,
 533 Lai Hong Wong, Pavol Drotár, Adam Kosiorek, Andrew Senior, Richard Tanburn, Taylor Ap-
 534 plebaum, Souradeep Basu, Demis Hassabis, and Pushmeet Kohli. AlphaGenome: advancing
 535 regulatory variant effect prediction with a unified DNA sequence model. URL <https://www.biorxiv.org/content/10.1101/2025.06.25.661532v2>. ISSN: 2692-8205 Pages:
 536 2025.06.25.661532 Section: New Results.
- 537 Yu Bo, Weian Mao, Yanjun Shao, Weiqiang Bai, Peng Ye, Xinzhu Ma, Junbo Zhao, Hao Chen, and
 538 Chunhua Shen. Revisiting convolution architecture in the realm of DNA foundation models. URL

- 540 <http://arxiv.org/abs/2502.18538>.
541
- 542 Sebastian Carrasco Pro, Heather Hook, David Bray, Daniel Berenzy, Devlin Moyer, Meimei
543 Yin, Adam Thomas Labadorf, Ryan Tewhey, Trevor Siggers, and Juan Ignacio Fuxman Bass.
544 Widespread perturbation of ETS factor binding sites in cancer. 14(1):913. ISSN 2041-
545 1723. doi: 10.1038/s41467-023-36535-8. URL [https://www.nature.com/articles/
546 s41467-023-36535-8](https://www.nature.com/articles/s41467-023-36535-8). Publisher: Nature Publishing Group.
- 547 Hannah Carter, Sining Chen, Leyla Isik, Svitlana Tyekucheva, Victor E. Velculescu, Kenneth W.
548 Kinzler, Bert Vogelstein, and Rachel Karchin. Cancer-specific high-throughput annotation of
549 somatic mutations: Computational prediction of driver missense mutations. 69(16):6660–6667.
550 ISSN 0008-5472. doi: 10.1158/0008-5472.CAN-09-1133. URL [https://doi.org/10.
551 1158/0008-5472.CAN-09-1133](https://doi.org/10.1158/0008-5472.CAN-09-1133).
- 552 Andrea Castro, Kivilcim Ozturk, Rachel Marty Pyke, Su Xian, Maurizio Zanetti, and Hannah Carter.
553 Elevated neoantigen levels in tumors with somatic mutations in the HLA-a, HLA-b, HLA-c and
554 b2m genes. 12(6):107. ISSN 1755-8794. doi: 10.1186/s12920-019-0544-1. URL <https://doi.org/10.1186/s12920-019-0544-1>.
- 555
556 Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. URL
557 <http://arxiv.org/abs/2307.08691>.
558
- 559 Giovanna Dashi and Markku Varjosalo. Oncofusions – shaping cancer care. 30(1):oyae126. ISSN
560 1549-490X. doi: 10.1093/oncolo/oyae126. URL [https://doi.org/10.1093/oncolo/
561 oyaee126](https://doi.org/10.1093/oncolo/oyae126).
- 562 Shiron Drusinsky, Sean Whalen, and Katherine S. Pollard. Deep-learning prediction of gene expres-
563 sion from personal genomes. URL [https://www.biorxiv.org/content/10.1101/
564 2024.07.27.605449v1](https://www.biorxiv.org/content/10.1101/2024.07.27.605449v1). Pages: 2024.07.27.605449 Section: New Results.
- 565
566 Prashant Gupta, Aashi Jindal, Gaurav Ahuja, Jayadeva, and Debarka Sengupta. A new deep learning
567 technique reveals the exclusive functional contributions of individual cancer mutations. 298(8):
568 102177. ISSN 0021-9258. doi: 10.1016/j.jbc.2022.102177. URL [https://www.ncbi.nlm.
569 nih.gov/pmc/articles/PMC9304782/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9304782/).
- 570 Douglas Hanahan. Hallmarks of cancer: New dimensions. 12(1):31–46. ISSN 2159-8274.
571 doi: 10.1158/2159-8290.CD-21-1059. URL [https://doi.org/10.1158/2159-8290.
572 CD-21-1059](https://doi.org/10.1158/2159-8290.CD-21-1059).
- 573 Johannes C Hingerl, Alexander Karollus, and Julien Gagneur. Flashzoi: An enhanced borzoi for ac-
574 celerated genomic analysis. pp. btaf467. ISSN 1367-4811. doi: 10.1093/bioinformatics/btaf467.
575 URL <https://doi.org/10.1093/bioinformatics/btaf467>.
- 576
577 Mengying Hu and Maria Chikina. InstaPrism: an r package for fast implementation of BayesPrism.
578 40(7):btae440. ISSN 1367-4811. doi: 10.1093/bioinformatics/btae440. URL [https://doi.
579 org/10.1093/bioinformatics/btae440](https://doi.org/10.1093/bioinformatics/btae440).
- 580 Connie Huang, Richard W. Shuai, Parth Baokar, Ryan Chung, Ruchir Rastogi, Pooja Kathail, and
581 Nilah M. Ioannidis. Personal transcriptome variation is poorly explained by current genomic
582 deep learning models. 55(12):2056–2059. ISSN 1546-1718. doi: 10.1038/s41588-023-01574-w.
583 URL <https://www.nature.com/articles/s41588-023-01574-w>. Publisher: Na-
584 ture Publishing Group.
- 585 Riku Katainen, Kashyap Dave, Esa Pitkänen, Kimmo Palin, Teemu Kivioja, Niko Välimäki, Alexan-
586 dra E. Gylfe, Heikki Ristolainen, Ulrika A. Hänninen, Tatiana Cajuso, Johanna Kondelin, Tomas
587 Tanskanen, Jukka-Pekka Mecklin, Heikki Järvinen, Laura Renkonen-Sinisalo, Anna Lepistö, Eevi
588 Kaasinen, Outi Kilpivaara, Sari Tuupanen, Martin Enge, Jussi Taipale, and Lauri A. Aaltonen.
589 CTCF/cohesin-binding sites are frequently mutated in cancer. 47(7):818–821. ISSN 1546-1718.
590 doi: 10.1038/ng.3335. URL <https://www.nature.com/articles/ng.3335>. Pub-
591 lisher: Nature Publishing Group.
- 592 Sunkyu Kim, Heewon Lee, Keonwoo Kim, and Jaewoo Kang. Mut2vec: distributed representation
593 of cancerous mutations. 11(2):33. ISSN 1755-8794. doi: 10.1186/s12920-018-0349-7. URL
<https://doi.org/10.1186/s12920-018-0349-7>.

- 594 Qiyuan Li, Ji-Heui Seo, Barbara Stranger, Aaron McKenna, Itsik Pe'er, Thomas LaFramboise,
595 Myles Brown, Svitlana Tyekucheva, and Matthew L. Freedman. Integrative eQTL-based anal-
596 yses reveal the biology of breast cancer risk loci. 152(3):633–641. ISSN 0092-8674, 1097-
597 4172. doi: 10.1016/j.cell.2012.12.034. URL [https://www.cell.com/cell/abstract/
598 S0092-8674\(12\)01556-5](https://www.cell.com/cell/abstract/S0092-8674(12)01556-5). Publisher: Elsevier.
- 599
600 Johannes Linder, Divyanshi Srivastava, Han Yuan, Vikram Agarwal, and David R. Kelley. Predicting
601 RNA-seq coverage from DNA sequence as a unifying model of gene regulation. URL [https://
602 //www.biorxiv.org/content/10.1101/2023.08.30.555582v1](https://www.biorxiv.org/content/10.1101/2023.08.30.555582v1).
- 603 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, a. URL [http://
604 arxiv.org/abs/1711.05101](http://arxiv.org/abs/1711.05101).
- 605
606 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts, b. URL
607 <http://arxiv.org/abs/1608.03983>.
- 608
609 Daria Ostroverkhova, Teresa M. Przytycka, and Anna R. Panchenko. Cancer driver mutations: pre-
610 dictions and reality. 29(7):554–566. ISSN 1471-4914, 1471-499X. doi: 10.1016/j.molmed.
611 2023.03.007. URL [https://www.cell.com/trends/molecular-medicine/
612 abstract/S1471-4914\(23\)00067-9](https://www.cell.com/trends/molecular-medicine/abstract/S1471-4914(23)00067-9). Publisher: Elsevier.
- 613
614 Joel S. Parker, Michael Mullins, Maggie C.U. Cheang, Samuel Leung, David Voduc, Tammi Vick-
615 ery, Sherri Davies, Christiane Fauron, Xiaping He, Zhiyuan Hu, John F. Quackenbush, Inge J.
616 Stijleman, Juan Palazzo, J.S. Marron, Andrew B. Nobel, Elaine Mardis, Torsten O. Nielsen,
617 Matthew J. Ellis, Charles M. Perou, and Philip S. Bernard. Supervised risk predictor of breast can-
618 cer based on intrinsic subtypes. 27(8):1160–1167. ISSN 0732-183X. doi: 10.1200/JCO.2008.18.
619 1370. URL <https://ascopubs.org/doi/10.1200/JCO.2008.18.1370>. Publisher:
Wolters Kluwer.
- 620
621 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
622 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-
623 ward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
624 Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance
625 deep learning library. URL <http://arxiv.org/abs/1912.01703>.
- 626
627 Gili Perry, Maya Dadiani, Smadar Kahana-Edwin, Anya Pavlovski, Barak Markus, Gil Hor-
628 nung, Nora Balint-Lahat, Ady Yosepovich, Goni Hout-Siloni, Jasmine Jacob-Hirsch, Miri
629 Sklair-Levy, Eitan Friedman, Iris Barshack, Bella Kaufman, Einav Nili Gal-Yam, and
630 Shani Paluch-Shimon. Divergence of mutational signatures in association with breast can-
631 cer subtype. 61(11):1056–1070. ISSN 1098-2744. doi: 10.1002/mc.23461. URL
632 <https://onlinelibrary.wiley.com/doi/abs/10.1002/mc.23461>. eprint:
<https://onlinelibrary.wiley.com/doi/pdf/10.1002/mc.23461>.
- 633
634 Allison Piovesan, Francesca Antonaros, Lorenza Vitale, Pierluigi Strippoli, Maria Chiara Pelleri,
635 and Maria Caracausi. Human protein-coding genes and gene feature statistics in 2019. 12(1):315.
636 ISSN 1756-0500. doi: 10.1186/s13104-019-4343-8. URL [https://doi.org/10.1186/
637 s13104-019-4343-8](https://doi.org/10.1186/s13104-019-4343-8).
- 638
639 Gustav Alexander Poulsen, Simon Grund Sørensen, Randi Istrup Juul, Morten Muhlig Nielsen,
640 and Jakob Skou Pedersen. Sequence dependencies and mutation rates of localized mutational
641 processes in cancer. 15(1):63. ISSN 1756-994X. doi: 10.1186/s13073-023-01217-z. URL
<https://doi.org/10.1186/s13073-023-01217-z>.
- 642
643 Peter Priestley, Jonathan Baber, Martijn P. Lolkema, Neeltje Steeghs, Ewart de Bruijn, Charles
644 Shale, Korneel Duyvesteyn, Susan Haidari, Arne van Hoeck, Wendy Onstenk, Paul Roepman,
645 Mircea Voda, Haiko J. Bloemendal, Vivianne C. G. Tjan-Heijnen, Carla M. L. van Herpen, Mari-
646 ette Labots, Petronella O. Witteveen, Egbert F. Smit, Stefan Sleijfer, Emile E. Voest, and Edwin
647 Cuppen. Pan-cancer whole-genome analyses of metastatic solid tumours. 575(7781):210–216.
ISSN 1476-4687. doi: 10.1038/s41586-019-1689-y. URL [https://www.nature.com/
articles/s41586-019-1689-y](https://www.nature.com/articles/s41586-019-1689-y). Publisher: Nature Publishing Group.

- 648 Michael Pudjihartono, Nicholas Pudjihartono, Justin M. O’Sullivan, and William Schierding.
649 Melanoma-specific mutation hotspots in distal, non-coding, promoter-interacting regions
650 implicate novel candidate driver genes. 131(10):1644–1655. ISSN 1532-1827.
651 doi: 10.1038/s41416-024-02870-w. URL [https://www.nature.com/articles/
652 s41416-024-02870-w](https://www.nature.com/articles/s41416-024-02870-w). Publisher: Nature Publishing Group.
- 653 Ruchir Rastogi, Aniketh Janardhan Reddy, Ryan Chung, and Nilah M. Ioannidis. Fine-
654 tuning sequence-to-expression models on personal genome and transcriptome data. URL
655 <https://www.biorxiv.org/content/10.1101/2024.09.23.614632v1>. Pages:
656 2024.09.23.614632 Section: New Results.
- 657 Prima Sanjaya, Katri Maljanen, Riku Katainen, Sebastian M. Waszak, J. C. Ambrose, P. Aru-
658 mugam, R. Bevers, M. Bleda, F. Boardman-Pretty, C. R. Boustred, H. Brittain, M. A. Brown,
659 M. J. Caulfield, G. C. Chan, A. Giess, J. N. Griffin, A. Hamblin, S. Henderson, T. J. P. Hub-
660 bard, R. Jackson, L. J. Jones, D. Kasperaviciute, M. Kayikci, A. Kousathanas, L. Lahnstein,
661 A. Lakey, S. E. A. Leigh, I. U. S. Leong, F. J. Leong, F. Maleady-Crowe, M. McEntagart,
662 F. Minneci, J. Mitchell, L. Moutsianas, M. Mueller, N. Murugaesu, A. C. Need, P. O’Donovan,
663 C. A. Odhams, C. Patch, D. Perez-Gil, M. B. Perez-Gil, J. Pullinger, T. Rahim, A. Rendon,
664 T. Rogers, K. Savage, K. Sawant, R. H. Scott, A. Siddiq, A. Siddiq, S. C. Smith, A. Sosinsky,
665 A. Stuckey, M. Tanguy, A. L. Taylor Tavares, E. R. A. Thomas, S. R. Thompson, A. Tucci,
666 M. J. Welland, E. Williams, K. Witkowska, S. M. Wood, M. Zarowiecki, Lauri A. Aaltonen,
667 Oliver Stegle, Jan O. Korbel, Esa Pitkänen, and Genomics England Research Consortium.
668 Mutation-attention (MuAt): deep representation learning of somatic mutations for tumour typ-
669 ing and subtyping. 15(1):47. ISSN 1756-994X. doi: 10.1186/s13073-023-01204-4. URL
670 <https://doi.org/10.1186/s13073-023-01204-4>.
- 671 Alexander Sasse, Bernard Ng, Anna E. Spiro, Shinya Tasaki, David A. Bennett, Christopher Gaiteri,
672 Philip L. De Jager, Maria Chikina, and Sara Mostafavi. Benchmarking of deep neural networks for
673 predicting personal gene expression from DNA sequence highlights shortcomings. 55(12):2060–
674 2064. ISSN 1546-1718. doi: 10.1038/s41588-023-01524-6. URL [https://www.nature.
675 com/articles/s41588-023-01524-6](https://www.nature.com/articles/s41588-023-01524-6). Publisher: Nature Publishing Group.
- 676 Anthony R. Soltis, Nicholas W. Bateman, Jianfang Liu, Trinh Nguyen, Teri J. Franks, Xijun Zhang,
677 Clifton L. Dalgard, Coralie Viollet, Stella Somiari, Chunhua Yan, Karen Zeman, William J.
678 Skinner, Jerry S. H. Lee, Harvey B. Pollard, Clesson Turner, Emanuel F. Petricoin, Daoud
679 Meerzaman, Thomas P. Conrads, Hai Hu, Rebecca Blackwell, Gauthaman Sukumar, Dagmar
680 Bacikova, Camille Alba, Elisa McGrath, Sraavya Poliseti, Meila Tuck, Alden Chiu, Gabe
681 Peterson, Caroline Larson, Leonid Kvecher, Brenda Deyarmin, Jennifer Kane, Katie Miller,
682 Kelly A. Conrads, Brian L. Hood, Sasha C. Makohon-Moore, Tamara S. Abulez, Elisa Baldelli,
683 Mariaelena Pierobon, Qing-rong Chen, Henry Rodriguez, Sean E. Hanlon, Anthony R. Soltis,
684 Nicholas W. Bateman, Jianfang Liu, Trinh Nguyen, Teri J. Franks, Xijun Zhang, Clifton L.
685 Dalgard, Coralie Viollet, Stella Somiari, Chunhua Yan, Karen Zeman, William J. Skinner,
686 Jerry S. H. Lee, Harvey B. Pollard, Clesson Turner, Emanuel F. Petricoin, Daoud Meerzaman,
687 Thomas P. Conrads, Hai Hu, Craig D. Shriver, Christopher A. Moskaluk, Robert F. Browning,
688 Matthew D. Wilkerson, Craig D. Shriver, Christopher A. Moskaluk, Robert F. Browning, and
689 Matthew D. Wilkerson. Proteogenomic analysis of lung adenocarcinoma reveals tumor hetero-
690 geneity, survival determinants, and therapeutically relevant pathways. 3(11):100819. ISSN 2666-
691 3791. doi: 10.1016/j.xcrm.2022.100819. URL [https://www.sciencedirect.com/
692 science/article/pii/S2666379122003780](https://www.sciencedirect.com/science/article/pii/S2666379122003780).
- 693 Alona Sosinsky, John Ambrose, William Cross, Clare Turnbull, Shirley Henderson, Louise Jones,
694 Angela Hamblin, Prabhu Arumugam, Georgia Chan, Daniel Chubb, Boris Noyvert, Jonathan
695 Mitchell, Susan Walker, Katy Bowman, Dorota Pasko, Marianna Buongiorno Pereira, Nadezda
696 Volkova, Antonio Rueda-Martin, Daniel Perez-Gil, Javier Lopez, John Pullinger, Afshan Sid-
697 diq, Tala Zainy, Tasnim Choudhury, Olena Yavorska, Tom Fowler, David Bentley, Clare Kings-
698 ley, Sandra Hing, Zandra Deans, Augusto Rendon, Sue Hill, Mark Caulfield, and Nirupa Mu-
699 rugaesu. Insights for precision oncology from the integration of genomic and clinical data of
700 13,880 tumors from the 100,000 genomes cancer programme. 30(1):279–289. ISSN 1546-
701 170X. doi: 10.1038/s41591-023-02682-0. URL [https://www.nature.com/articles/
s41591-023-02682-0](https://www.nature.com/articles/s41591-023-02682-0). Publisher: Nature Publishing Group.

- 702 Anna E. Spiro, Xinming Tu, Yilun Sheng, Alexander Sasse, Rezwane Hosseini, Maria Chikina,
703 and Sara Mostafavi. A scalable approach to investigating sequence-to-function predictions from
704 personal genomes. URL [https://www.biorxiv.org/content/10.1101/2025.02.](https://www.biorxiv.org/content/10.1101/2025.02.21.639494v3)
705 [21.639494v3](https://www.biorxiv.org/content/10.1101/2025.02.21.639494v3). ISSN: 2692-8205 Pages: 2025.02.21.639494 Section: New Results.
706
- 707 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. RoFormer: En-
708 hanced transformer with rotary position embedding. URL [http://arxiv.org/abs/2104.](http://arxiv.org/abs/2104.09864)
709 [09864](http://arxiv.org/abs/2104.09864).
- 710 Lakshman Sundaram, Arvind Kumar, Matthew Zatzman, Adriana Salcedo, Neal Ravindra, Shadi
711 Shams, Bryan H. Louie, S. Tansu Bagdatli, Matthew A. Myers, Shahab Sarmashghi, Hyo Young
712 Choi, Won-Young Choi, Kathryn E. Yost, Yanding Zhao, Jeffrey M. Granja, Toshinori Hinoue,
713 D. Neil Hayes, Andrew Cherniack, Ina Felau, Hani Choudhry, Jean C. Zenklusen, Kyle Kai-How
714 Farh, Andrew McPherson, Christina Curtis, Peter W. Laird, The Cancer Genome Atlas Analy-
715 sis Network, M. Ryan Corces, Howard Y. Chang, and William J. Greenleaf. Single-cell chro-
716 matin accessibility reveals malignant regulatory programs in primary human cancers. 385(6713):
717 eadk9217. doi: 10.1126/science.adk9217. URL [https://www.science.org/doi/10.](https://www.science.org/doi/10.1126/science.adk9217)
718 [1126/science.adk9217](https://www.science.org/doi/10.1126/science.adk9217).
- 719 Adrian Tan, Gonçalo R. Abecasis, and Hyun Min Kang. Unified representation of genetic variants.
720 31(13):2202–2204. ISSN 1367-4803. doi: 10.1093/bioinformatics/btv112. URL [https://](https://doi.org/10.1093/bioinformatics/btv112)
721 doi.org/10.1093/bioinformatics/btv112.
- 722 Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu. Diagnosis of
723 multiple cancer types by shrunken centroids of gene expression. 99(10):6567–6572. doi:
724 10.1073/pnas.082099299. URL [https://www.pnas.org/doi/full/10.1073/pnas.](https://www.pnas.org/doi/full/10.1073/pnas.082099299)
725 [082099299](https://www.pnas.org/doi/full/10.1073/pnas.082099299). Publisher: Proceedings of the National Academy of Sciences.
726
- 727 Sjoerd J. D. Tjalsma, Niels J. Rinzema, Marjon J. A. M. Verstegen, Michelle J. Robers, Andrea
728 Nieto-Aliseda, Richard A. Gremmen, Amin Allahyar, Mauro J. Muraro, Peter H. L. Krijger,
729 and Wouter de Laat. Long-range enhancer-controlled genes are hypersensitive to regulatory fac-
730 tor perturbations. 5(3). ISSN 2666-979X. doi: 10.1016/j.xgen.2025.100778. URL [https:](https://www.cell.com/cell-genomics/abstract/S2666-979X(25)00034-5)
731 [//www.cell.com/cell-genomics/abstract/S2666-979X\(25\)00034-5](https://www.cell.com/cell-genomics/abstract/S2666-979X(25)00034-5). Pub-
732 lisher: Elsevier.
- 733 John N. Weinstein, Eric A. Collisson, Gordon B. Mills, Kenna R. Mills Shaw, Brad A. Ozenberger,
734 Kyle Ellrott, Ilya Shmulevich, Chris Sander, and Joshua M. Stuart. The cancer genome atlas
735 pan-cancer analysis project. 45(10):1113–1120. ISSN 1546-1718. doi: 10.1038/ng.2764. URL
736 <https://www.nature.com/articles/ng.2764>. Number: 10 Publisher: Nature Pub-
737 lishing Group.
- 738 Ross Wightman. `fastai/timmdocs`. URL <https://github.com/fastai/timmdocs>.
739 `original-date: 2021-01-19T18:56:46Z`.
740
- 741 Chunfu Xiao, Xiaoge Liu, Peiyu Liu, Xinwei Xu, Chao Yao, Chunqiong Li, Qi Xiao, Tiannan
742 Guo, Li Zhang, Yongjun Qian, Chao Wang, Yiting Dong, Yingxuan Wang, Zhi Peng, Chuanhui
743 Han, Qiang Cheng, Ni A. An, and Chuan-Yun Li. Oncogenic roles of young human *de novo*
744 genes and their potential as neoantigens in cancer immunotherapy. 5(9):100928. ISSN 2666-
745 979X. doi: 10.1016/j.xgen.2025.100928. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S2666979X25001843)
746 [science/article/pii/S2666979X25001843](https://www.sciencedirect.com/science/article/pii/S2666979X25001843).
- 747 Wei Zhang, Ana Bojorquez-Gomez, Daniel Ortiz Velez, Guorong Xu, Kyle S. Sanchez, John Paul
748 Shen, Kevin Chen, Katherine Licon, Collin Melton, Katrina M. Olson, Michael Ku Yu, Justin K.
749 Huang, Hannah Carter, Emma K. Farley, Michael Snyder, Stephanie I. Fraley, Jason F. Kreisberg,
750 and Trey Ideker. A global transcriptional network connecting noncoding mutations to changes
751 in tumor gene expression. 50(4):613–620. ISSN 1546-1718. doi: 10.1038/s41588-018-0091-2.
752 URL <https://www.nature.com/articles/s41588-018-0091-2>. Publisher: Na-
753 ture Publishing Group.
- 754 Dixian Zhu, Tianbao Yang, and Livnat Jerby. Gradient aligned regression via pairwise losses. URL
755 <http://arxiv.org/abs/2402.06104>.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A APPENDIX

A.1 HYPERPARAMETERS

To train GVF, we conducted a minimal, manual hyperparameter search using the 5 training folds. We ultimately chose to use AdamW (Loshchilov & Hutter, a) with a learning rate of 10^{-4} , β s of (0.9, 0.95), and ϵ of 10^{-8} . We also used a cosine learning rate scheduler (Loshchilov & Hutter, b) from TIMM (Wightman) starting and ending the learning rate at 10^{-5} with 1 warm up epoch and complete decay after 3 epochs, with no cycles. The final GVF model had 4 transformer layers and roughly 1M parameters.

A.2 LLM USAGE

We used LLMs for drafting, polishing language, and literature search.