

Article

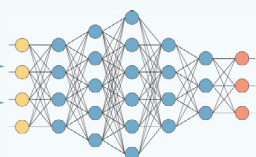
Deep Learning Implicitly Handles Tissue Specific Phenomena to Predict Tumor DNA Accessibility and Immune Activity

Predict DNA accessibility

Tumor RNA-seq profile

...TAATCTA... ...CACACC... ...GGAATCA...

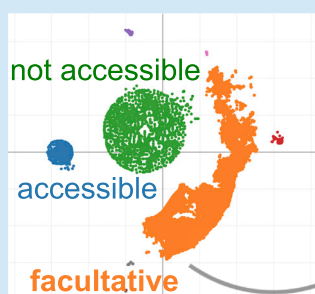
Loci DNA sequences



DNA accessibility

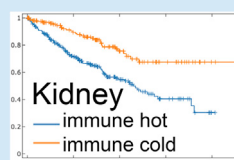
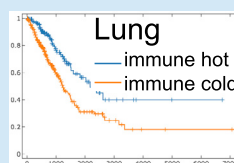
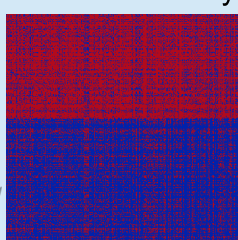
...TAATCTA... ...CACACC... ...GGAATCA...

DNA accessibility predictions stratify patients into immune-related subtypes



Tumor Sample

Loci Accessibility



Kamil Wnuk,
Jeremi Sudol,
Kevin B.
Givechian, Patrick
Soon-Shiong,
Shahrooz
Rabizadeh,
Christopher Szeto,
Charles Vaske

kamil.wnuk@immunitybio.com

HIGHLIGHTS

Tissue-specific aspects of DNA can be predicted in new tissue types if given RNA-seq

DNA accessibility prediction is most reliable at promoter and promoter flank regions

Clustering global chromatin state highlights immune pathway activity in tumors

Accessibility patterns discriminate immune active tumors that can differ in prognosis

Wnuk et al., iScience 20, 119–136
October 25, 2019 © 2019 The Author(s).
<https://doi.org/10.1016/j.isci.2019.09.018>

Article

Deep Learning Implicitly Handles Tissue Specific Phenomena to Predict Tumor DNA Accessibility and Immune Activity

Kamil Wnuk,^{1,4,*} Jeremi Sudol,¹ Kevin B. Givechian,² Patrick Soon-Shiong,¹ Shahrooz Rabizadeh,¹ Christopher Szeto,³ and Charles Vaske³

SUMMARY

DNA accessibility is a key dynamic feature of chromatin regulation that can potentiate transcriptional events and tumor progression. To gain insight into chromatin state across existing tumor data, we improved neural network models for predicting accessibility from DNA sequence and extended them to incorporate a global set of RNA sequencing gene expression inputs. Our expression-informed model expanded the application domain beyond specific tissue types to tissues not present in training and achieved consistently high accuracy in predicting DNA accessibility at promoter and promoter flank regions. We then leveraged our new tool by analyzing the DNA accessibility landscape of promoters across The Cancer Genome Atlas. We show that in lung adenocarcinoma the accessibility perspective uniquely highlights immune pathways inversely correlated with a more open chromatin state and that accessibility patterns learned from even a single tumor type can discriminate immune inflammation across many cancers, often with direct relation to patient prognosis.

INTRODUCTION

DNA accessibility plays a key role in the regulatory machinery of DNA transcription. Locations where DNA is not tightly bound in nucleosomes, detectable as DNase I hypersensitivity sites (DHSs), render the sequence accessible to other DNA-binding proteins, including a wide range of transcription factors (TFs). DHSs are cell specific and play a crucial role in determining transcriptional events that differentiate cells.

Furthermore, genome-wide association studies (GWAS) have revealed that the vast majority of genetic variants significantly associated with many diseases and traits are located in non-coding regions (Deplancke et al., 2016) and well over half non-coding single nucleotide polymorphisms (SNPs) affect DHSs (Maurano et al., 2012). Thus variable access to DNA regulatory elements plays a key role not only in normal cell development but also in altered expression profiles associated with disease states (Deplancke et al., 2016; Xu et al., 2015), including cancer.

In an effort to go beyond association studies and gain deeper insight into how changes in DNA sequence affect transcriptional regulation, some groups have developed predictive models for a multitude of genomic phenomena. Several works have recently made significant advances in accuracy of such DNA-sequence-based prediction tasks by applying neural network models to problems such as TF binding (Alipanahi et al., 2015; Zhou and Troyanskaya, 2015; Quang and Xie, 2016; Lanchantin et al., 2016), promoter-enhancer interactions (Singh et al., 2016), DNA accessibility (Hoffman et al., 2018; Kelley et al., 2016; Zhou and Troyanskaya, 2015), and DNA methylation states (Angermueller et al., 2017; Hoffman et al., 2018).

One common issue that limits the broad applicability of these models is the cell-type-specific nature of many of the underlying biological mechanisms, such as DHSs. All the above examples approached their prediction problems by learning to estimate the conditional probability, $p(a|d, b)$, where a is the accessibility (or other attributes) of a segment of DNA sequence, d , and b is a discrete label of tissue type. In practice, this meant either training a separate model for each cell or tissue type or having a single model output multiple tissue-specific (multi-task) predictions. This made it difficult to apply the models to new data and limits them from being integrated into broader scope pathway models (Vaske et al., 2010).

Conveniently, a number of studies have demonstrated that gene expression levels from RNA sequencing (RNA-seq) can be used to discriminate cell types (Sudmant et al., 2015; Breschi et al., 2016; Conesa et al.,

¹ImmunityBio Inc., Culver City, CA 90232, USA

²NantOmics LLC, Culver City, CA 90232, USA

³ImmunityBio Inc., Santa Cruz, CA 95060, USA

⁴Lead Contact

*Correspondence:
kamil.wnuk@immunitybio.com

<https://doi.org/10.1016/j.isci.2019.09.018>



2016), providing evidence that $p(b|r)$ can be learned (where r is a vector of RNA-seq gene expression measurements). In addition, DNase sequencing (DNase-seq) and microarray-based gene expression levels from matched samples were found to cluster similarly according to biological relationships, with many DHSs found to significantly correlate with gene expressions (Sheffield et al., 2013).

Our work focuses on overcoming the barrier to broad applicability due to cell-type-specific phenomena by putting the burden on a deep neural network classifier to handle the complex relationship between expression and DNA sequence accessibility without intermediate discrete tissue labels. Our model directly estimates $p(a|d, r)$, and thus implicitly handles the space of possible tissue types and states, B , because $p(a|d, r) = \sum_{i \in B} p(a|d, b_i) p(b_i|r)$. This allows the model to exploit similarity information in the space of tissue types and make predictions for previously unseen tissues whose gene expressions were similar but unique from samples in the training data.

We build on the Basset neural network architecture of Kelley et al., which recently demonstrated state-of-the-art results on DNA accessibility prediction (Kelley et al., 2016). They factored the cell-specific DHS issue into their work by first creating a binary matrix of sample tissues and their respective accessibility for a list of genomic sites. The universal list of (potentially accessible) sites was found by agglomerative clustering of all overlapping DNase-seq peaks across all samples before training. These potentially accessible sites defined the 600-base pair DNA segments used as inputs. Second, they set up the model's final layer as a multitask output, with a distinct prediction unit for each tissue type.

We began by showing that neural network performance on the original task of Kelley et al., predicting accessibility at held-out sites across 164 tissue types, could be improved by strategic factorization of convolutional layers. Then, based on our hypothesis that a neural network should be capable of learning to modulate appropriately its prediction of the DNase-seq signal if informed with a global RNA-seq state, we extended the network to handle a vector of gene expression values as input. To train this model, we constructed a new dataset of samples from the ENCODE project (ENCODE Project Consortium, 2012) consisting only of RNA-seq and DNase-seq measurements whose correspondence could be determined. Our new model achieved compelling results for held-out tissue types and proved to be very reliable at predicting accessibility at promoter and promoter flank genomic regions.

We then applied our accessibility prediction model to whole genomes across six cancer cohorts from The Cancer Genome Atlas (TCGA) (TCGA, 2018), as summarized in Figure 1A, and highlighted that accessibility complements RNA-seq. For example, clustering lung adenocarcinoma (LUAD) samples based on predicted accessibility was distinct from any RNA-seq-based cluster assignment and revealed a group of patients showing enrichment for pathways involved in immune response. Furthermore, splitting the same cohort by immune cell composition revealed a difference in survival and enabled training of a classifier to detect an immune-inflamed tumor state in LUAD, based only on accessibility predictions for a small set of sites. This classifier allowed us to discriminate immune-active patient groups with significant differences in survival in several distinct cancers, often aligning with findings from other cancer immunology studies. To the extent of our knowledge, this was the first time that a prediction of DNA accessibility had been applied to whole genomes in TCGA cancer cohorts to infer the chromatin landscape across cancers. In parallel with our work, accessibility was measured using assay for transposase-accessible chromatin using sequencing (ATAC-seq) for select samples across several TCGA cohorts (Corces et al., 2018), and we have shown that these new empirical results validate our predictions (Figure 7).

We anticipate that our expression-informed model not only may provide detailed information regarding DNA accessibility across tissues and enable discrimination of immune-inflamed tumors but also might be used to predict individual patient response to various immune-based therapies. We also expect our approach to be a useful tool in understanding other conditions where chromatin state is suspected to play an important role, including aging (Moskowitz et al., 2017), neurodegenerative disease (Berson et al., 2018), autoimmune diseases (Farh et al., 2014), as well as autism (Zhao et al., 2018). Finally, we stress that our approach to implicitly handling tissue type and state can be used in any DNA-based prediction task.

RESULTS

Convolutional Layer Factorization Improves Accuracy

As a baseline we first implemented our own version of the Basset architecture (Figure 1B), with minor changes not related to network structure. Compared with an receiver operating characteristic area under

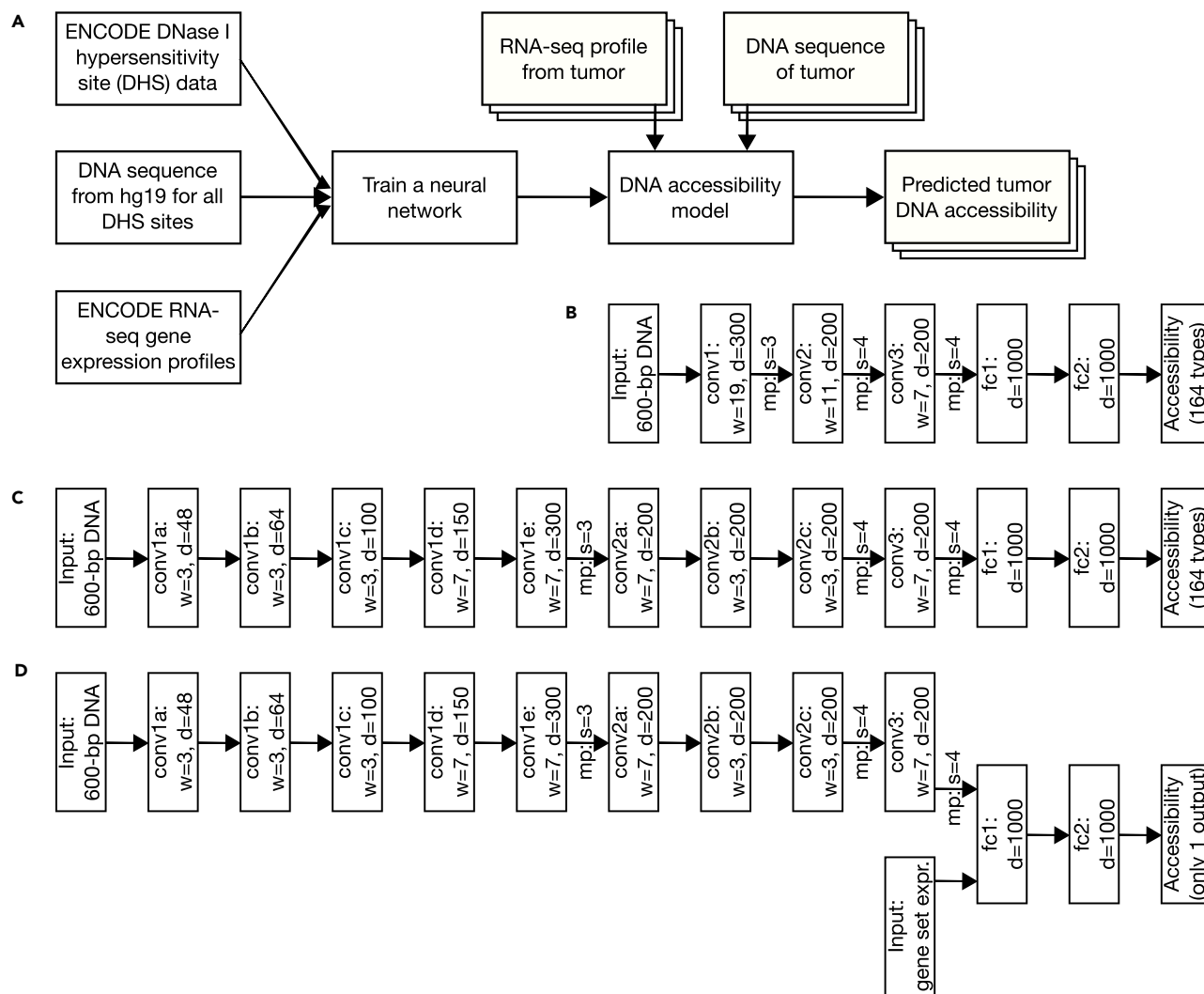


Figure 1. Overview of Our Pipeline from Training to Application

(A) DHS, hg19 DNA, and RNA-seq information are all used to train the neural network. With tumor RNA-seq and DNA-seq data input the DNA accessibility model can then be used to predict chromatin state in tumors.

(B–D) (B) The neural network architectures for the tissue-specific baseline model, (C) the tissue-specific factorized convolutions model, and (D) the expression-informed model are shown. Depth (d) is provided for all fully connected (fc) layers. Convolution (conv) layers also list their width (w). Max pooling (mp) is indicated where present between layers and is always applied with equal size and stride (s).

See also [Figures S1, S3, and S4](#).

the curve (ROC AUC) = 0.895 reported by Kelley et al. on their benchmark test set (Kelley et al., 2016), and confirmed by applying the pre-trained model provided by the authors, our baseline implementation achieved a mean ROC AUC = 0.903 (Table 1).

Following the success of many works demonstrating that deeper hierarchies of small convolutional kernels tend to improve neural network performance (Simonyan and Zisserman, 2014; Szegedy et al., 2016; He et al., 2016), we experimented with factorization of large convolutional layers in the baseline model. We found that factorization of layers closest to the data was the most significant for improving accuracy (Figure S1). When only the second convolutional layer was factorized, the speed of learning improved during the early epochs of training, but final accuracy was not noticeably affected compared with our baseline implementation. An overall improvement in both rate of learning and final accuracy was achieved when both the first and second convolutional layers were factorized (Figure 1C), leading to a mean ROC AUC = 0.910 (Table 1).

Tissue-Specific Model	Mean ROC AUC	Mean PR AUC
Basset (pre-trained)	0.895	0.561
Our baseline	0.903	0.582
Our factorized convolutions	0.910	0.605

Table 1. Tissue-Specific Model Results on Basset Benchmark Dataset

Furthermore, despite following the same training procedure and taking no additional steps to account for class imbalance, our final tissue-specific model's mean precision-recall area under the curve (PR AUC) = 0.605 compared favorably with the mean PR AUC = 0.561 reported as the best result obtained by the Basset model.

ENCODE DNase-Seq and RNA-Seq Dataset

To train a model for predicting accessibility that is informed implicitly about tissue state through gene expression it was necessary to build a new dataset where both DNase-seq and RNA-seq were available for a large and diverse collection of different tissue types. We collected all human samples from the ENCODE project (ENCODE Project Consortium, 2012) for which correspondence between RNA-seq and DNase-seq measurements could be determined. After errors were filtered out, the final dataset consisted of 74 unique tissue types, with 220 DNase-seq files and 304 RNA-seq files. A validation set of 5% randomly held-out samples was split from the data so that tissue types were diverse but still independent measurements of tissues also appeared in training.

Two sets of paired test and training partitions were created (Figure S2). The first partition pair (tissue overlap in Table 2) was constructed in the same way as the validation data by randomly holding out test samples and allowing for tissue type overlap with training. The second partition pair (held-out tissue in Table 2) was constructed such that the test set was composed only of samples from tissue types that were not present in either training or validation partitions. The latter was meant to more accurately simulate the intended application scenario and was thus the main focus throughout our analysis.

The data partitions were revised once from their first iteration when several erroneous samples were discovered and revoked by the ENCODE consortium. Once revoked samples were removed, we saw a significant decrease in spurious DHSs. Table 2 shows the distributions of the final dataset.

Expression-Informed Model Can Predict Accessibility in Held-Out Tissues

We explored several alternative versions of our expression-informed neural network along with a range of different hyperparameters. Based on validation performance, we found that concatenating the global gene set expression vector directly with output from the convolutional layers, using a large batch size with appropriately matched learning rate, and weight initialization from a model trained using more noisy training data made the most impact. Changing the fraction of positive samples per training batch from 0.5 to 0.25 also led to a minor improvement.

Table S2 shows that final model (Figure 1D) performance on the validation set, both overall and by tissue type, was consistent across each of the two training partitions with respect to both ROC AUC as well as PR AUC. Tables 3 and 4 summarize the results of applying our model across whole genomes, at all potential DHSs. For tissue types with more than a single file pair in the test set, each sample's results are listed.

As expected, overall the model was less accurate on completely new tissue types; however, even in the more challenging scenario, the overall PR AUC was higher than the best tissue-specific models evaluated on known tissue types. Note that several of the results in Table 4 were within similar ranges as predictions whose sample types overlapped with training.

Expression-Informed Model Predictions Are Highly Reliable at Promoter and Promoter Flank Genomic Sites

To better understand the performance characteristics and limitations of our model, we broke down our ENCODE validation and test results by genomic site type. Table 5 details the distribution of annotations

Partition	Unique Tissues	DNase-Seq Files	RNA-Seq Files
Validation	10	11	12
Tissue overlap train	73	195	277
Tissue overlap test	12	14	15
Held-out tissue train	66	198	281
Held-out tissue test	8	11	11

Table 2. File- and Tissue-Type Distribution per Dataset Partition

applied to the 1.71 million sites considered in the held-out tissue training set, the percentage of all positive samples that fall within each annotation, and the percentage of samples per each annotation type that are positive. Note that a single site may overlap with more than one annotation.

We found that even for samples in which the model performed poorly overall (Figure 2C), predictions within promoter and promoter flank regions consistently attained a high level of accuracy (Figures 2A and 2D, Table 6), achieving a PR AUC = 0.839 over all held-out tissue types and a PR AUC = 0.911 over randomly held-out samples (validation set).

We also confirmed that the accuracy of these predictions was independent of whether the promoter and promoter flank sites overlapped with the regions of genes used in our RNA-seq input gene set (Figure 2B). Selecting a threshold for classification of only promoter and promoter flank sites such that precision is 80% (20% false discovery rate) on the held-out tissue test set, our trained model recalls 65.3% of accessible promoter regions, with a false-positive rate of 10%. Applying this same threshold to the validation set where tissues are allowed to overlap with the training set, the model achieves a precision of 93.4%, recalling 62.6% of accessible promoter regions, and has a false-positive rate of only 3.5%.

We also investigated the accuracy at enhancer sites, finding a PR AUC = 0.732 over held-out tissues and PR AUC = 0.889 over randomly held-out samples (validation set). Differently from promoter and promoter flank regions, however, enhancer prediction accuracies showed a high variance between test samples (Figure 2D, Table S3). Thus, more investigation is necessary before relying on accessibility predictions at enhancers in further analysis.

To quantify the effect of similarity to training data on prediction performance we looked at correlation between PR AUC (computed independently for all predictions in each whole genome sample) and distance of each test and validation sample to its closest sample in the training set (Table 7). As might be expected from Figure 2D, we confirm that prediction performance is less correlated with test sample similarity to training data at promoter and promoter flank sites than when PR AUC is evaluated over all potentially accessible sites.

Promoter Accessibility Patterns across Cohorts from The Cancer Genome Atlas

We applied our trained model to promoter and promoter flank sites in TCGA samples from six cohorts (LUAD, LUSC, KICH, KIRC, KIRP, and BRCA). Across all samples in these cohorts for which whole-genome sequencing (WGS) was available, 3,172 interest regions had one SNP, 78 had two SNPs, and only 9 regions had between three and five SNPs. A total of 465 sites included insertion or deletions (INDELs), and only 7 sites featured both an INDEL and an SNP. Lung cancers exhibited the highest average number of mutated sites per patient from our selected cohorts (Figure 3A).

To observe the effect of region changes on accessibility, we compared predictions with and without SNPs and INDELs present. INDELs had the greatest impact on predicted accessibility, exhibiting a higher variance than SNPs (Figure 3B) and leading to a change in accessibility classification in 5.46% cases (at the previously defined accessibility threshold that achieved 80% precision) (Figure 3C). As there were generally few somatic mutations impacting accessibility prediction sites, and the percentage of those that actually impacted classifier decisions was even lower, we can conclude that any patterns we observed at the landscape scale of chromatin state predictions will be dictated more by gene expression levels providing global context of tissue state rather than somatic mutations.

Sample Tissue Type	ROC AUC	PR AUC
A172	0.959	0.721
Left renal pelvis	0.967	0.838
Small intestine	0.951, 0.926	0.737, 0.571
Muscle of arm	0.968	0.843
Forelimb muscle	0.968	0.808
Keratinocyte	0.939	0.644
Skin fibroblast	0.948, 0.947	0.770, 0.763
Large intestine	0.964	0.727
Muscle of back	0.967, 0.954	0.853, 0.854
Adrenal gland	0.957	0.743
SK-N-DZ	0.898	0.652
Fibroblast of lung	0.942	0.840
Mean tissue type AUC	0.950	0.758
Overall AUC	0.947	0.748

Table 3. Tissue Overlap Whole-Genome Test Results, with Scores Computed across All Tissues in Bold

To get a landscape view of how promoter and promoter flank sites behave, we embedded their binary accessibility decisions across all our selected TCGA samples in two dimensions with t-distributed stochastic neighbor embedding (t-SNE) (van der Maaten and Hintin, 2008). This clearly separated sites into constitutively accessible, constitutively not accessible, and facultative (Figure 3D). A few very small clusters of several hundred sites were also formed, which were groups of typically constitutive sites that acted uniquely in one or two individual patients.

Second, we stacked all predictions into a single vector per patient to form accessibility profiles for all samples in our six TCGA cohorts and again applied t-SNE to visualize relationships (Figure 3F). This qualitatively showed that looking at cancers from the viewpoint of DNA accessibility highlights different relationships than analysis of RNA-seq alone. For example, in the RNA-seq t-SNE space (Figure 3E) a clear separation emerges among breast cancers (BRCA), which correspond to basal-like versus luminal A/B and HER2-enriched clusters. In contrast, in accessibility t-SNE space (Figure 3F), the lung (LUAD, LUSC) and breast (BRCA) cancer samples appear to share some common characteristics. There also appears to be a slight partition into left and right groups in how lung cancer samples arrange in the embedding space, with LUAD forming a distinct subset away from the LUSC/BRCA modality within the lung/breast supercluster. Some similar clustering trends were observed in empirical ATAC-seq measurements performed concurrent to our work (Corces et al., 2018).

One of many potential biological factors that may contribute to overlap in the accessibility space is the impact of hormone activity on both breast and lung tumors; this activity in turn is epigenetically regulated (Zhang and Ho, 2011), thus some chromatin state patterns could be shared. Cell type composition is another factor that is likely to play a role in determining our observations of both expression and accessibility states of tissues; however, at this time we cannot isolate which biological confounders that determine tissue state contribute in what quantities to accessibility predictions.

Accessibility Is Linked to Immune Activity in Lung Adenocarcinoma

We subsequently explored the biological associations of our model's accessibility predictions by examining transcriptomic data from LUAD samples, as this tumor type has been shown to be of particular interest in chromatin accessibility studies due to the impact on progression (Polak et al., 2014; Kim and Kim, 2016). Upon clustering LUAD samples for which WGS was available according to their predicted accessibility, clear bifurcation into low- (C0, 21 samples) and high-accessibility (C1, 20 samples) samples was

Sample Tissue Type	ROC AUC	PR AUC
Left kidney	0.965	0.778
OCI-LY7	0.899, 0.899, 0.886, 0.886	0.654, 0.654, 0.655, 0.654
Prostate gland	0.865	0.516
Hindlimb muscle	0.943	0.824
Spleen	0.913	0.582
Astrocyte	0.919, 0.944	0.787, 0.613
Fibroblast of skin of abdomen	0.964	0.826
G401	0.739, 0.846	0.459, 0.516
Mean tissue type AUC	0.898	0.655
Overall AUC	0.897	0.621

Table 4. Held-Out Tissue Whole-Genome Test Results, with Scores Computed across All Tissues in Bold

observed (Figure 4A). Cluster assignment based on predicted accessibility was distinct from any cluster assignments using the same methodology on RNA-seq directly (Figure 4B).

Differential KEGG pathway expression analysis with Enrichr (Kuleshov et al., 2016) showed the Chemokine Signaling Pathway (hsa04062) to be upregulated in the low DNA accessibility (C0) patient group. This association held true whether using TOIL RNA-seq data (Vivian et al., 2017; TOIL RNA-seq Recompute, 2016) (Enrichr adjusted p value [adj. p] = 1.191×10^{-6}) or HiSeqV2 RNA-seq data (TCGA Genome Characterization Center UNC, 2017) (Enrichr adj. p = 0.0145) (see methods and Figure S5 for details). Chemokines are involved in multiple key processes in tumor growth and immune response (Nagarsheth et al., 2017; Rivas-Fuentes et al., 2015; Sarvaiya et al., 2013), and their regulation by epigenetic mechanisms has also previously been reported (Flavahan et al., 2017; Yasmin et al., 2015).

No difference in tumor mutation burden was found between the two clusters (two-sided t test: $t = -0.696$, $p = 0.491$), but interestingly the C0 group exhibited higher expression of immune checkpoint genes (Figure 4C). Cell type enrichment analysis (Aran et al., 2017) of lymphoids and myeloids also revealed a higher level of class-switched memory B cells (two-sided t test: $t = 4.040$, $p = 0.000385$, Benjamini-Hochberg [BH] adj. $p = 0.0131$) in C0, although estimated levels were generally low in both clusters. Other immune cell estimates exhibited no differences with adj. $p < 0.1$, which was largely limited by small sample set size.

Total Number of Accessible Sites Correlates with Activity in Immune Pathways

To enhance the scope of our findings, we extended our analysis to all LUAD patient samples for which WGS was not available by predicting accessibility using just the reference genome (hg19/GRCh37) and gene expression data. Although no mutation information was included for these additional samples, only 37 of 5,449 sites (6.79%) used to cluster all WGS data included any instances of mutations. With the additional consideration that only a small percentage of all mutations actually flip binary class predictions (Figure 3C), it is unlikely that cluster assignment of new non-WGS samples was significantly affected by this missing information. As before, the expanded set of patient samples was clustered into two groups according to accessibility.

The group with generally lower accessibility (C0) again exhibited generally higher checkpoint levels (Figures 4D and 4E). However, visualizing the first three principal components and coloring points by total number of accessible promoter and promoter flank sites (Figure 4F), we did find a smooth change in value along a continuous manifold of samples, primarily along the first principal component (Spearman correlation = 0.989, $p = 0.0$).

Therefore, instead of differential analysis, all protein-coding genes were filtered by correlation with the total number of accessible sites and evaluated for KEGG pathway enrichment. We found that all genes satisfying the threshold (correlation absolute value >0.4) had negative correlation values. As some relationships

Site Type	% of All Sites	% of All Positive Examples	% per Site Type that Are Positive
Exon	3.47	5.08	9.74
Intragenic	49.94	45.25	6.04
Intergenic	39.89	34.11	5.70
Promoter and flank	6.37	30.33	31.75
Enhancer	1.08	3.81	23.47
Other	5.39	5.20	6.43

Table 5. Distribution of Potentially Accessible Sites by Annotation

may not be linear but still monotonic, we focus on the Spearman measure (Table S4), although the Pearson measure yielded similar top pathways (Table S5). Osteoclast Differentiation (hsa04380, adj. $p = 7.45 \times 10^{-15}$) was the most significantly correlated pathway.

Interestingly, the process of osteoclast differentiation is controlled by two essential cytokines (Kim and Kim, 2016): macrophage colony stimulating factor and the receptor activator of nuclear factor (NF)- κ B ligand. Tumor Necrosis Factor (TNF) Signaling Pathway (hsa04668, adj. $p = 1.30 \times 10^{-8}$) also appeared among the top three results. TNF has a pro-inflammatory effect and has been noted to play a critical role in the control of apoptosis, angiogenesis, proliferation, invasion, and metastasis (Yasmin et al., 2015).

When pathways were sorted by significance, the Chemokine Signaling Pathway, observed in WGS-only cluster analysis, appeared 11th (adj. $p = 4.01 \times 10^{-5}$). Other notable pathways appearing in the top 10 included Pathways in Cancer (hsa05200, adj. $p = 6.50 \times 10^{-7}$), Regulation of Actin Cytoskeleton (hsa04810, adj. $p = 6.50 \times 10^{-7}$), NF- κ B Signaling Pathway (hsa04064, adj. $p = 2.68 \times 10^{-7}$), and Epstein-Barr Virus Infection (hsa05169, adj. $p = 2.83 \times 10^{-5}$). Focal Adhesion (hsa04510, adj. p value = 2.29×10^{-4}) appeared 20th but is worth noting as it resurfaces in later analysis.

Majority of Genes with Differential Accessibility Exhibit Consistent Differential Expression in Immune Cell-Driven Clusters

To investigate accessibility patterns specifically in the context of different tumor immune environments, all LUAD samples were clustered into two groups according to lymphoid and myeloid levels based on xCell cell type enrichment analysis (Aran et al., 2017). Lymphoid and myeloid cells were selected for their roles in the adaptive and innate immune system, respectively. A thin margin was introduced between clusters to exclude samples with near-ambiguous label assignment (Figure S6).

Patients in X0 (141 samples) were enriched for many immune cells (Figure 5A), as well as checkpoint gene expression (Figure 5C), and tended to have narrower distributions of both number of accessible promoter and flank sites (generally lower than X1, two-sided t test $p = 1.07 \times 10^{-3}$) and total overall methylation (generally higher than X1, two-sided t test $p = 1.29 \times 10^{-7}$) (Figure 5B). These samples also reflected significantly favorable survival (Figure 5G). We therefore interpreted X0 as the group of “immune-hot” patients in LUAD and X1 as “immune-cold” patients.

After eliminating sites that exhibited low standard deviation, we selected all significantly differentiated accessibility sites between the two clusters and mapped them to their nearest gene. Qualitatively, we observed that several groups of sites act together in different ways and that those different clusters of chromatin state behavior are stable across significance thresholds (Figures 5D and 5E). We found that when a majority of sites corresponding to a single gene were accessible more frequently in one cluster, that gene exhibited upregulated expression in the same cluster most of the time (64.7% for genes more accessible in X0 and 64.2% in X1).

Genes whose expression was consistent with increased accessibility in X0 showed near-significant levels of enrichment for some pathways that had previously surfaced in our correlation results such as Focal Adhesion (hsa04510, adj. $p = 0.0355$) and Osteoclast Differentiation (hsa04380, adj. $p = 0.0936$) (Table S6). No

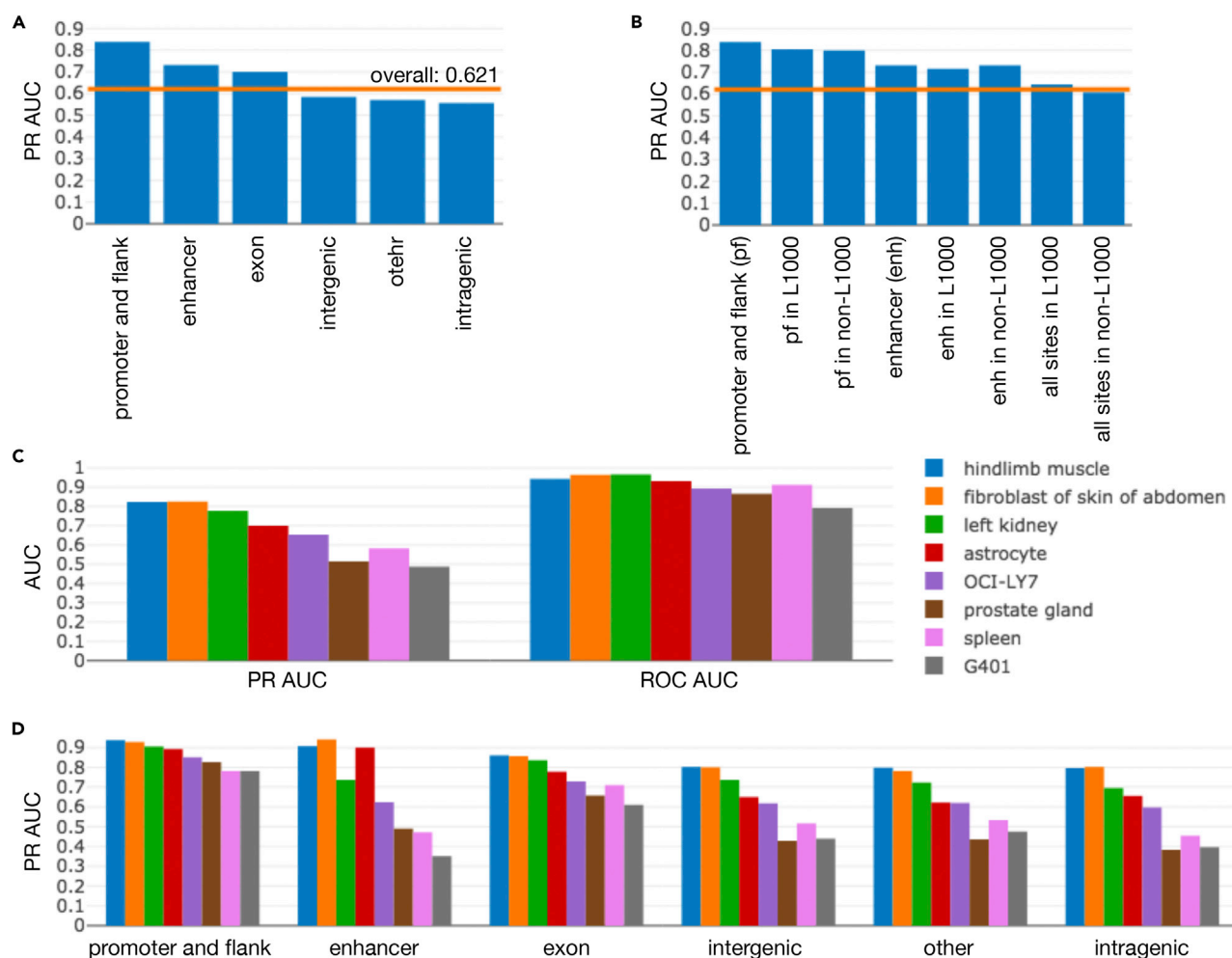


Figure 2. Promoter and Promoter Flank Accessibility Is Highly Predictable, but Enhancers Show Variability

(A) Promoter flank (pf) accessibility is highly predictable (PR AUC = 0.839), as shown by the genomic site performance breakdown over all samples in the held-out tissues test set. The orange line indicates overall PR AUC computed across all test samples and all sites.

(B) No clear performance difference was observed when genomic sites across the held-out tissue test set were split into those that did (in L1000) and did not (non-L1000) overlap the L1000 RNA-seq input gene set. Note that not all sites overlapped with known gene regions, so the union of the L1000 and non-L1000 subsets did not always make up the complete set of sites of a certain type.

(C) Overall metrics separated by tissue type show that some held-out tissues in the test set were more challenging as reflected by lower AUCs.

(D) Predictions at enhancers were highly variable between samples, even with good PR AUC, and performance on pf regions remained consistently high, even for tissues where overall results were lowest.

significant pathways were found for genes consistent with increased accessibility in X1, or those inconsistent with more accessibility in X0. The strongest significance in pathway enrichment existed in the set of genes inconsistent with increased accessibility in X1 (up in X0 despite accessibility predictions voting for upregulation in X1) (Table S7). The most prominent of the enriched pathways in this group were Platelet Activation (hsa04611, adj. $p = 4.38 \times 10^{-4}$), Inflammatory Mediator Regulation of TRP Channels (hsa0475, adj. $p = 0.0109$), and several with adj. $p = 0.0235$: Chemokine Signaling Pathway (hsa04062), Focal Adhesion (hsa04510), cGMP-PKG Signaling Pathway (hsa04022), Intestinal Immune Network for IgA Production (hsa04672), and Vascular Smooth Muscle Contraction (hsa04270).

These findings suggest that in LUAD tumors, partial regulation of immune- and cytokine-controlled pathways may be exerted via an activator mechanism at promoters. Furthermore, a more significant component of chemokine signaling and platelet activation that distinguishes immune-active patients may be subject to repressor regulation at promoter sites.

Sample Tissue Type	ROC AUC	PR AUC
Left kidney	0.949	0.905
OCI-LY7	0.869, 0.868, 0.864, 0.864	0.842, 0.842, 0.859, 0.859
Prostate gland	0.897	0.826
Hindlimb muscle	0.935	0.938
Spleen	0.867	0.782
Astrocyte	0.925, 0.914	0.946, 0.838
Fibroblast of skin of abdomen	0.951	0.929
G401	0.798, 0.828	0.807, 0.757
Overall AUC	0.876	0.839

Table 6. Held-Out Tissue Test Results Restricted to Promoter and Promoter Flank Sites, with Scores Computed across All Tissues in Bold

Patterns of Promoter Accessibility Predict Immune-Hot Tumors with Impact on Patient Survival across Several Cancers

To further explore the link between DNA accessibility, immune activity, and clinical outcomes, we trained an ensemble of three classifiers to detect an immune-hot tumor state in LUAD based only on a small subset of accessibility predictions. Applying the ensemble to all of LUAD (Figure 6A) led to a cleaner and more significant partition of patients (compared to Figures 5F and S6F) into immune-hot and immune-cold tumors. Further applying the classifier ensemble to accessibility predictions across 11 other cancers in TCGA revealed cases wherein the immune-hot state learned from LUAD was beneficial to patient survival (Figures 6A–6E), detrimental to survival (Figures 6J and 6K), or had little impact (Figures 6F–6I and 6L), with varying degrees of significance.

Our findings aligned very well with a comprehensive analysis of immune subtypes across TCGA (Thorsson et al., 2018), which characterized the influence of immune activations on survival. Despite very different methodologies, their plots also indicated that in cohorts such as LUAD, SKCM, and CESC activation of immune subtypes was associated with better outcomes; that the opposite was true in KIRC, LGG, and STAD; and that little impact was visible in LUSC. They did not, however, discuss accessibility as a potential additional biomarker for immune activity.

Significant negative impact of immune activity on survival in KIRC was also shown in a separate cohort of clinical data from Oulu University Hospital (Mella et al., 2015), confirming that this trend is not unique to TCGA. Interestingly, the study used CD8⁺ T cell count cutoffs to stratify patients with renal cell carcinoma into two groups. Based on xCell estimates this was the second most significantly enriched immune cell type (two-sided t test BH adj. $p = 3.57 \times 10^{-26}$) in KIRC immune-hot patients identified by our classifier, after activated dendritic cells (two-sided t test BH adj. $p = 1.07 \times 10^{-26}$) (Figure S8). Although the training cohort (LUAD) did express some difference in CD8⁺ T cell enrichment scores between hot and cold tumors (two-sided t test BH adj. $p = 9.65 \times 10^{-15}$), it was not in the top 10 most significantly different immune cells (Figure 5A). From this we see that our classifiers operating on accessibility predictions learned a more complex decision boundary than simply focusing on direct correlates with the most differentiated immune cell compositions in the training set.

An additional curiosity specific to the KIRC partition from our classifier ensemble is that little difference in CD274 (also called PD-L1) and PDCD1LG2 (also called PD-L2) expression is visible between the two predicted classes; however, the strong difference in PDCD1 (also called PD-1) expression levels that exists in other immune-hot versus immune-cold partitions does exist. This unique state of checkpoint-related gene expression may be linked to the low response rate in patients with renal cell carcinoma to anti-PD-L1 therapies, compared with more favorable responses found for anti-PD-1 drugs, in early-phase clinical trials (Weinstock and McDermott, 2015). Interestingly, an empirical study of ATAC-seq peaks across cancers (Corces et al., 2018) also found a link between four regulatory regions that exhibited distinct chromatin accessibility patterns across cancers and expression of CD274.

PR AUC Evaluation Domain	Pearson Correlation	Pearson p Value	Spearman ρ	Spearman p Value
Overall	-0.7472	1.77×10^{-5}	-0.7080	7.52×10^{-5}
Promoter and flank	-0.6795	1.87×10^{-4}	-0.5417	5.16×10^{-3}

Table 7. Correlation of PR AUC with Test Sample Distance to the Nearest Training Sample

ATAC-Seq Measurements Validate Predictions in TCGA Samples

In parallel to our analysis, the chromatin state of select samples across several TCGA cohorts was empirically measured using ATAC-seq (Corces et al., 2018). At a minimum overlap threshold of 70% we found that 10.9% (61,342 of 562,709) of all pan-cancer peaks identified in the study corresponded directly with 56.3% of our (108,970) promoter and promoter flank sites at which we applied our model. At lower peak overlap thresholds this percentage increased significantly; for example, at 10% minimum overlap 83.6% of our promoter and promoter flank sites had corresponding pan-cancer ATAC-seq peaks. For correspondences defined by the 70% overlap threshold we first showed that the means of normalized count values for individual ATAC-seq peaks across lung and kidney cohorts had clearly distinct distributions (Figure 7A) between constitutive and facultative site categories (identified based on clustering our TCGA predictions, Figure 3D). Constitutive sites had consistently high mean peak counts in the accessible category and consistently low mean peak counts in the not accessible category, whereas facultative sites corresponded to ATAC-seq peaks with a broader distribution of mean normalized count values centered between the previous two categories (Figure 7A).

We then explored the distributions of ATAC-seq peak counts as stratified directly by our accessibility classifier predictions in addition to the above site categories (Figures 7B and 7C). In this case no means were computed across samples; every prediction site in every TCGA sample was considered as one data point. Normalized counts for peaks corresponding to sites predicted as accessible were generally distributed at higher values than peaks corresponding to sites predicted as not accessible, and the difference between predicted accessibility distributions was significantly more striking at constitutive sites than at facultative sites. This observation held true for all cohorts used in our immune classification experiments (Figure 6) for which there existed TCGA samples with both ATAC-seq measurements and our accessibility predictions. This result suggests that chromatin sites whose accessibility changes dynamically within a tissue type tend to remain closer to an accessibility threshold compared with less dynamic sites. The distributions did differ some in shape between cohorts, which may partly be explained by the fact that our constitutive and facultative site category labels used for this experiment were derived only based on predictions in lung and kidney cohorts; thus in different tissues wherein the dynamics of chromatin slightly differ those particular labels may not be as representative.

DISCUSSION

We have demonstrated that predictive models operating on DNA sequence data, additionally conditioned on a global set of RNA-seq gene expression inputs, can predict DHSs in unseen tissue types in a way that allows application to new samples without re-training. We showed that these models were capable of achieving consistently high performance for predictions at promoter and promoter flank regions of the genome. Leveraging this new tool for analysis of tumor genomes across different cell and tissue types, we provided a unique perspective on the DNA accessibility landscape across TCGA data. Complementary to the exploration of the full range of variable accessibility sites across cancers made possible by empirical measurements (Corces et al., 2018), our analysis of sites at reliably predictable genomic regions explored a more limited and thus more subtle set of chromatin dynamics, which proved to still be very information rich. Despite the more limited view of accessibility sites in our case, both studies found some similar clustering trends and both concluded that chromatin state plays a significant role in cancer immune response.

DNA accessibility is one of many factors that determine expression, which makes inversion of the relationship not trivial; knowing expression levels does not uniquely define the pattern of DHSs. Our expression-informed model (Figure 1D) learns a most likely mechanism by which the DNA sequence immediately

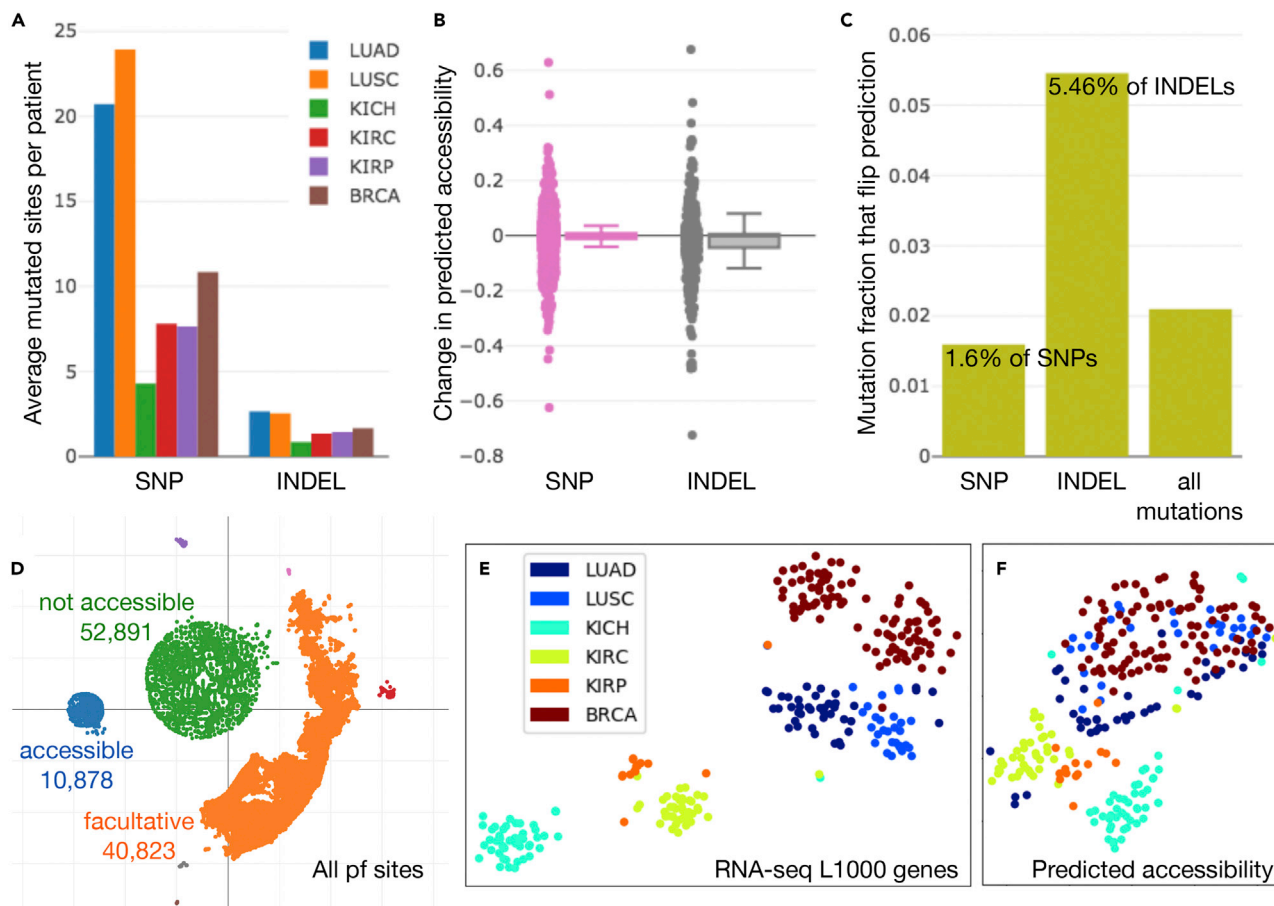


Figure 3. SNP and INDEL Mutations and Predicted Accessibility Landscape in Tumors

(A) The average number of SNP and insertion or deletion (INDEL) mutations that overlap prediction sites per patient across six TCGA cohorts is shown. (B) When predictions at sites with mutations were compared with and without applying mutations to the input DNA sequence, the change in predicted accessibility exhibited a higher variance for INDELS than SNPs.

(C) In addition, a larger fraction of sites with INDELS were responsible for a change in the classification decision (flipped prediction) than the fraction of sites with SNPs.

(D) Using t-SNE (perplexity = 50) to visualize the predicted accessibility of individual promoter flank (pf) sites across our selected TCGA samples, we identified which sites were facultative (orange), constitutively accessible (blue), and constitutively not accessible (green).

(E and F) (E) Finally, t-SNE applied to patient samples exhibited different relationships (such as a clear split in BRCA samples) when based on RNA-seq gene expression of the L1000 gene set, than (F) when based on predicted accessibility at all pf sites within each sample (in which case lung and breast cancers appeared to share some common characteristics).

surrounding a potential DHS determines its accessibility, conditioned also on an observed global expression state. Therefore, accessibility prediction applied across the whole genome is an approach to approximately invert gene expression to obtain most likely DHSs.

Our results showed that viewing tumors by promoter accessibility highlights immune pathways that would otherwise be harder to detect from completely unsupervised analysis of RNA-seq data alone. For example, we found several pathways inversely correlated with an overall more open chromatin state. Through identification of facultative accessibility sites linked with differential gene expression in immune-inflamed LUAD tumors and training of a classifier ensemble, we showed that patterns of predicted chromatin state at a small subset of genomic regions are predictive of immune activity across many tumor types, with direct implications for patient prognosis. We see such predictive models playing a significant future role in matching patients to appropriate immunotherapy treatment regimens, as well as in analysis of other conditions wherein epigenetic state may play a significant role, such as autoimmune disease, autism, aging, and neurodegenerative disease.

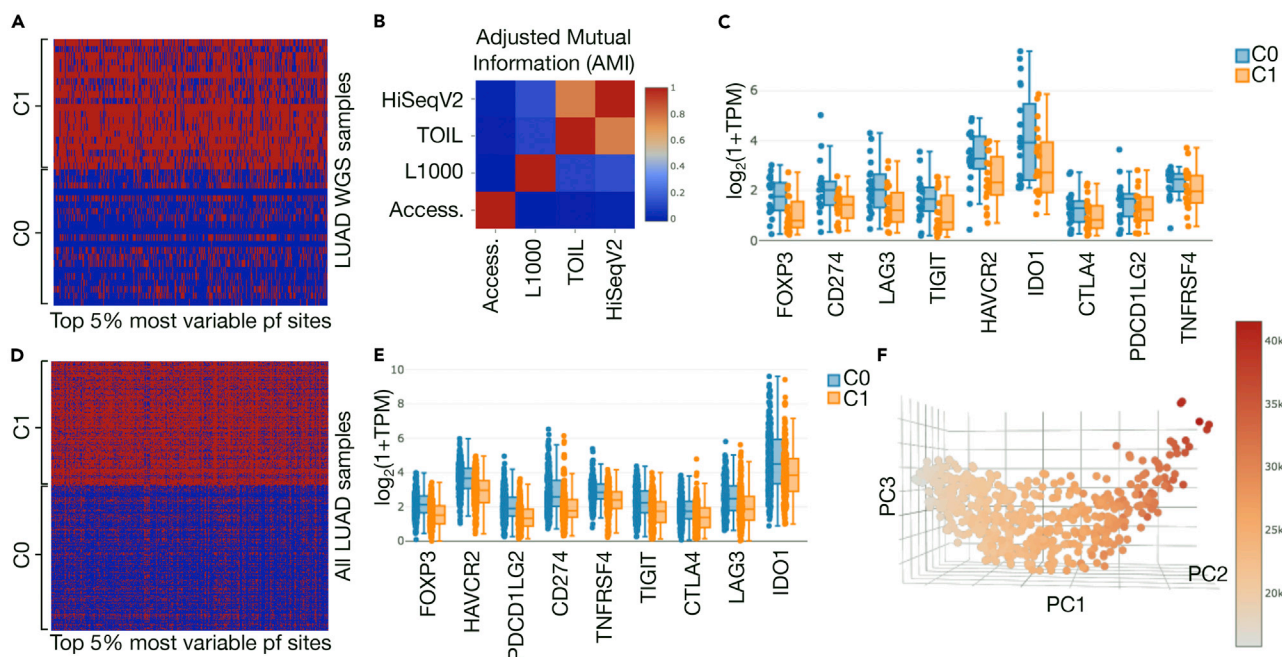


Figure 4. Promoter and Promoter Flank Accessibility and Checkpoint Gene Expression in LUAD WGS Samples Only and Augmented with Non-WGS Samples

(A) The heatmap and patient sample cluster assignment based on the top 5% most variable promoter and promoter flank (pf) accessibility sites across LUAD samples with WGS available are shown. Cluster 0 (C0) has lower overall accessibility (blue = not accessible), and cluster 1 (C1) exhibits generally higher accessibility (red = accessible).

(B) Adjusted mutual information (AMI) (1) between label assignments based on different data shows higher values (red) between different RNA-seq cluster assignments and low values (blue) between accessibility (Access.) and clusters based on any other data type.

(C) Distribution of key checkpoint gene expression levels (with x axis sorted by significance of two-sided t test between C0 and C1) shows that the low-accessibility group tends to have higher checkpoint levels.

(D) Applying the same procedure to the full LUAD cohort, which also includes predictions for all non-WGS samples, we see a similar split into low- (C0) and high (C1)-accessibility groups.

(E) The same trend in checkpoint expression is observed, with FOXP3 again appearing as the most significant difference (two-sided t test with Benjamini-Hochberg adjusted $p = 4.53 \times 10^{-19}$).

(F) Plotting promoter and flank accessibility with respect to its first three principal components (PC1–3) and coloring points by total number of accessible sites in a sample reveals a smoothly varying relationship, motivating a correlation-based approach to exploring the relationship between overall accessibility and gene expression levels.

It may also be interesting to pursue a deeper functional investigation of genes linked with accessibility. Genes with consistent behavior to accessibility are candidates that may be regulated via an activator mechanism at promoters, whereas genes with inconsistent behavior may be subject to alternative gene repression mechanisms, e.g., silencer elements or suppression via microRNAs.

In a few TCGA cohorts, our ensemble classification approach only identified a very small number of immune-hot tumor samples, making survival analysis impossible. The generalizability of our immune-related chromatin state across cancers was undoubtedly limited by only having trained the support vector machines (SVMs) on a single cohort, because immune cell composition and definition of an immune-active state varies across cancers (Thorsson et al., 2018); going forward, we will integrate accessibility signatures from multiple cohorts to train a more comprehensive subtyping of immune state.

Ideally, WGS for each of the samples in our training dataset should have been used to learn the most faithful representation of the true biology, as using only reference genome data introduces non-random noise in the input space. Unfortunately, such individual whole-genome data were not available for this project. Nonetheless, our work and that of others demonstrates that useful predictors can be learned despite this noise. Unlike models with multitask outputs, our architecture can easily support such individualized training without any changes, and when possible, it will be instituted in the future.

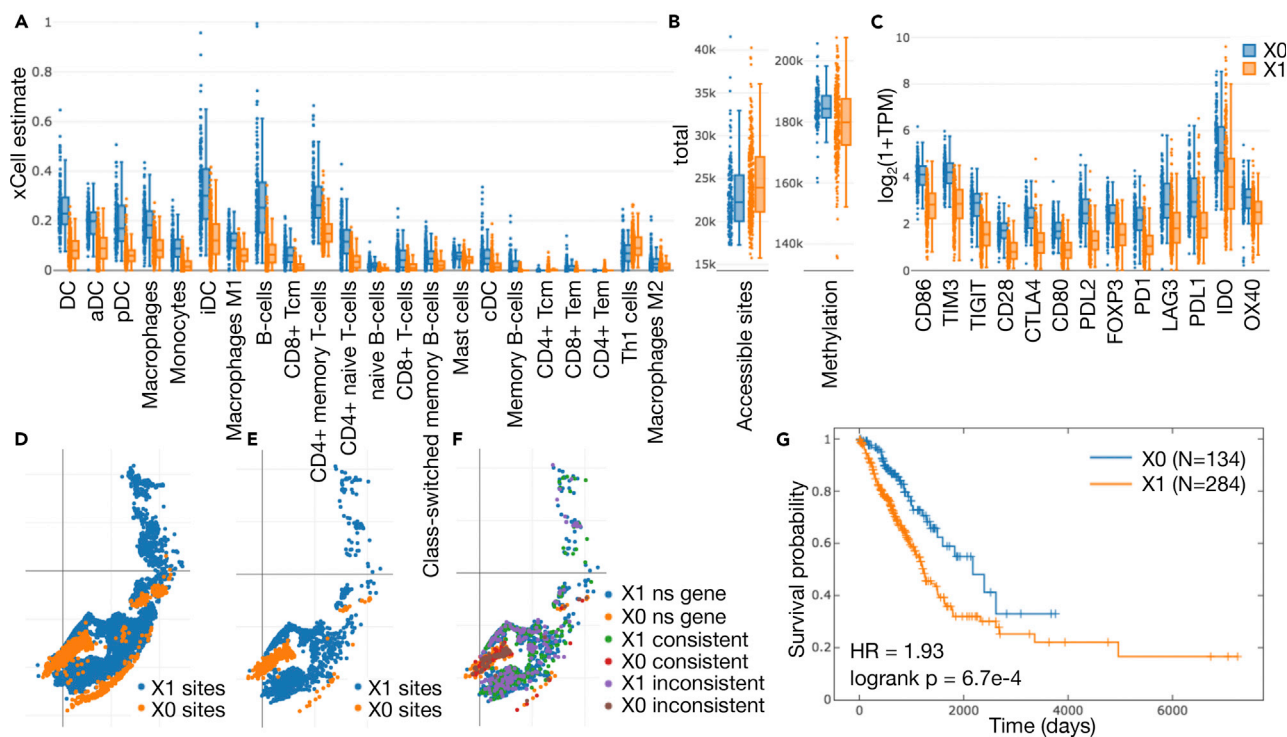


Figure 5. Enrichment in LUAD xCell-Derived Clusters (after Adding a Small Margin) by Cell Type, Checkpoint Expression, Methylation, Accessibility, and Survival

(A) Cell type enrichment distributions sorted by significance of two-sided t test for the two clusters (X0, X1), based on xCell lymphoid and myeloid cells, with Benjamini-Hochberg adjusted p value $< 1.0 \times 10^{-5}$ are shown.

(B) Total number of accessible promoter and promoter flank sites in each sample by cluster (two-sided t test $p = 1.07 \times 10^{-3}$) along with total methylation (two-sided t test $p = 1.29 \times 10^{-7}$).

(C) Checkpoint expression distributions, likewise sorted by significance, also point to a general difference in immune landscape between the two groups.

(D and E) (D) All sites with differences in accessibility based on a two-sided t test with Benjamini-Hochberg adjusted p values < 0.01 and (E) $< 1.0 \times 10^{-5}$ are illustrated on the t-SNE plot of promoter and promoter flank facultative sites. Sites with a difference satisfying the thresholds were assigned to the cluster in which they were more accessible.

(F) Accessibility differences are further broken down by how they align with direction of upregulation of corresponding nearby genes (ns gene, no significant difference in matching gene; consistent, direction of significant accessibility and gene expression differences are consistent; inconsistent, direction of significant accessibility and gene expression are inconsistent).

(G) Kaplan-Meier plots demonstrate better survival among X0 (immune hot) patients, shown with log rank test p value and hazard ratio (HR) based on a Cox proportional hazards (CoxPH) model regression using class assignment as the only explanatory variable.

See also [Figures S6](#) and [S7](#), and [Tables S4–S7](#).

We saw high variance for enhancer sites, but these sites are also interesting with respect to chromatin state and immunotherapy, because they have been linked with T cell dysfunction with potential for therapeutic reprogrammability in mice (Philip et al., 2017). At this time, it needs to be determined whether the large variance in performance is due to limitations in the model, noise in the data, or lack of necessary information in the available inputs. To this end, we look forward to future exploration of a more complete set of genes instead of a manually curated set, such as the LINCS L1000. Many alternatives exist to learn gene embeddings as part of model training, and we believe that ultimately an approach that efficiently incorporates all genes as input will be most effective.

Furthermore, there are a multitude of alternative model architectures such as residual connections (He et al., 2016), densely connected convolutional networks (Huang et al., 2016), and recurrent neural networks (Hochreiter and Schmidhuber, 1997) with additions such as attention (Bahdanau et al., 2014; Xu et al., 2015), which we believe are likely to improve performance of our model. These have been left for future evaluation, such as one rigorous study that has independently verified and built on our architectural innovations (Nair et al., 2019). The key contribution of this work was movement beyond the cell-type-specific limitations

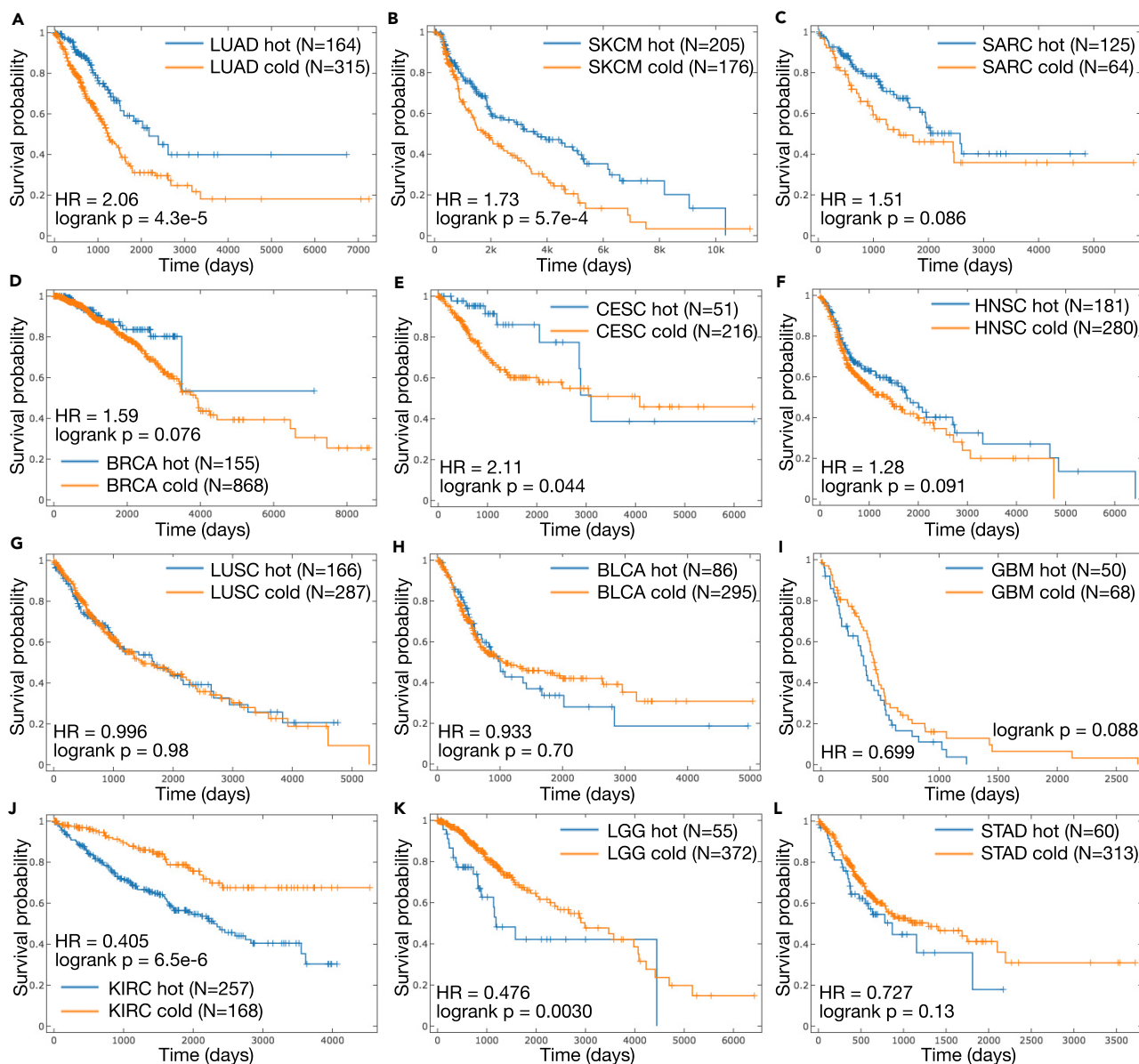


Figure 6. Application of the Three SVM Ensembles for Classification of Immune-Hot Tumors (Trained on Subsets of LUAD) with the Only Input Being a Vector of 484 Promoter and Flank Predicted Accessibility Decisions

All Kaplan-Meier plots show group size (N) for patients of both predicted immune activity classes (hot/cold) that satisfy a confidence threshold (see [Transparent Methods](#)). Also provided are log rank test p values and hazard ratio (HR) based on a Cox proportional hazards (CoxPH) model regression using class assignment as the only explanatory variable. Note that the time axis range on subplots varies by cohort and that the immune-hot state learned based on LUAD is not always beneficial for patient survival in other tumor types. Tumor types included (A) LUAD, lung adenocarcinoma; (B) SKCM, skin cutaneous melanoma; (C) SARC, sarcoma; (D) BRCA, breast invasive carcinoma; (E) CESC, cervical squamous cell carcinoma and endocervical adenocarcinoma; (F) HNSC, head and neck squamous cell carcinoma; (G) LUSC, lung squamous cell carcinoma; (H) BLCA, bladder urothelial carcinoma; (I) GBM, glioblastoma multiforme; (J) KIRC, kidney renal clear cell carcinoma; (K) LGG, brain lower grade glioma; and (L) STAD, stomach adenocarcinoma.

See also [Figure S8](#).

of DNA sequence classifiers, demonstration of the application of our expression-informed model to predict accessibility, and the ability of these predictions to distinguish prognostically alternative immune states across human cancers.

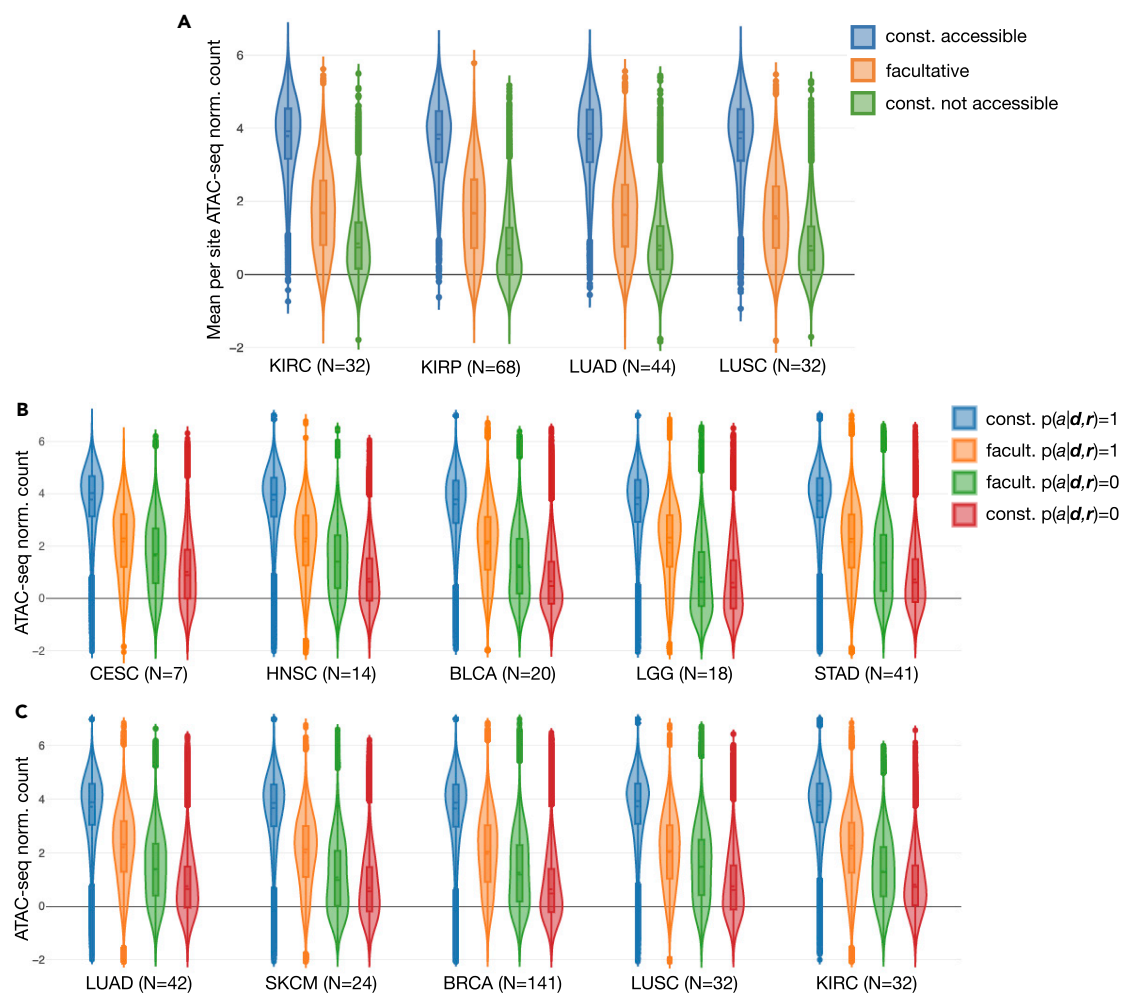


Figure 7. Validation of Our Promoter and Promoter Flank DNA Accessibility Predictions in TCGA with Empirical ATAC-Seq Measurements

(A) The top violin plots show the distributions of per ATAC-seq peak means of normalized counts in lung and kidney cohorts, for sites we labeled as constitutively (const.) accessible, facultative, or const. not accessible (based on our analysis shown in Figure 3D). Peak count values along all y axes were log transformed and quantile normalized as provided by the authors of the empirical study.

(B and C) Distributions of ATAC-seq peak normalized counts for all prediction sites across all available samples were further broken down per cohort by classification decision (accessible, $p(a|d,r) = 1$, and not accessible, $p(a|d,r) = 0$) in addition to site category. Site categories were either facultative (facult.) or constitutive (const.), the latter including both const. accessible as well as const. not accessible. The number of TCGA samples that contributed to each plot is shown (N =). (B and C) Only TCGA samples for which we had made predictions and were also empirically measured were used, but (A) utilized all available measured samples. The distribution plots were informed by $N * 61,342$ data points in (B and C), whereas for (A), where we considered the mean value for each site, there were only 61,342 data points total within each cohort.

Limitations of the Study

By design, convolutional neural networks only capture the influence of a small local neighborhood of DNA sequence (600 base pairs in our experiments) on predicted outputs, so impacts from more distal mutations on potential DNA accessibility sites can be captured implicitly only if they happen to influence expression levels of input of the global RNA-seq gene set. True biological function of DNA sequence may not be fully captured due to reliance on reference genome as a proxy for WGS in all training samples. At genomic regions where predictions demonstrate consistently good accuracy, we suspect that a fair amount of noise due to this approximation has averaged out over training data. Regions at which prediction accuracies have high variance across samples, such as enhancers, may be hard to predict for this reason, or other limitations with the training data or model assumptions. Both DNase-seq and RNA-seq measurements are taken from tissue samples, which feature heterogeneous cell type populations of varying proportions. The addition of RNA-seq data has enabled models to implicitly handle some degree of this type of noise;

however, due to the massive possible permutations in how such variations can manifest, trained models may not perform well in test cases where tissue types, compositions, or the RNA-seq expression quantification pipeline are drastically different than samples seen in training.

METHODS

All methods can be found in the accompanying [Transparent Methods supplemental file](#).

SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.isci.2019.09.018>.

ACKNOWLEDGMENTS

We would like to thank Patricia Spilman for editing support and Hermes Garban for discussions regarding the biology of chromatin and regulation at promoters.

AUTHOR CONTRIBUTIONS

K.W. conceived the project, curated the data, developed the model, performed analysis, and wrote the manuscript; J.S. contributed to data processing software; K.B.G. conceived and performed initial clustering analysis in LUAD and contributed to writing and editing; P.S.-S. secured funding; S.R. suggested RNA-seq as a suitable tissue-type signature; S.R., C.S., and C.V. provided cancer bioinformatics expertise and feedback and contributed to editing; C.V. and C.S. provided TCGA mutation calls.

DECLARATION OF INTERESTS

This work was funded by NantWorks affiliates (ImmunityBio, NantOmics, NantHealth) and performed by its employees; there are no other conflicts of interest.

Received: April 29, 2019

Revised: August 23, 2019

Accepted: September 11, 2019

Published: October 25, 2019

REFERENCES

- Alipanahi, B., Delong, A., Weirauch, M.T., and Frey, B.J. (2015). Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat. Biotechnol.* **33**, 831–838.
- Angermueller, C., Lee, H.J., Reik, W., and Stegle, O. (2017). DeepCpG: accurate prediction of single-cell DNA methylation states using deep learning. *Genome Biol.* **18**, 67.
- Aran, D., Hu, Z., and Butte, A.J. (2017). xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome Biol.* **18**, 220.
- Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arxiv.org*, arXiv:1409.0473.
- Berson, A., Nativio, R., Berger, S.L., and Bonini, N.M. (2018). Epigenetic regulation in neurodegenerative diseases. *Trends Neurosci.* **41**, 587–598.
- Breschi, A., Djebali, S., Gillis, J., Pervouchine, D.D., Dobin, A., Davis, C.A., Gingeras, T.R., and Guigo, R. (2016). Gene-specific patterns of expression variation across organs and species. *Genome Biol.* **17**, 151.
- Conesa, A., Madrigal, P., Tarazona, S., Gomez-Cabrero, D., Cervera, A., McPherson, A., Szczesniak, M.W., Gaffney, D.J., Elo, L.L., Zhang, X., and Mortazavi, A. (2016). A survey of best practices for RNA-seq data analysis. *Genome Biol.* **17**, 13.
- Corces, M.R., Granja, J.M., Shams, S., Louie, B.H., Seoane, J.A., Zhou, W., Silva, T.C., Groeneveld, C., Wong, C.K., Cho, S.W., et al. (2018). The chromatin accessibility landscape of primary human cancers. *Science* **362**, <https://doi.org/10.1126/science.aav1898>.
- Deplancke, B., Alpern, D., and Gardeux, V. (2016). The genetics of transcription factor DNA binding variation. *Cell* **166**, 538–554.
- ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74.
- Farh, K.K.-H., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W.J., Beik, S., Shores, N., Whitton, H., Ryan, R.J.H., Shishkin, A.A., et al. (2014). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* **518**, 337.
- He, K., Zhang, X., Ren, S. & Sun, J. 2016. Deep residual learning for image recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 770–778.
- Flavahan, W.A., Gaskell, E., and Bernstein, B.E. (2017). Epigenetic plasticity and the hallmarks of cancer. *Science* **357**, <https://doi.org/10.1126/science.aal2380>.
- Hochreiter, S., and Schmidhuber, J. (1997). Long short-term memory. *Neural Comput.* **9**, 1735–1780.
- Hoffman, G.E., Schadt, E.E., and Roussos, P. (2018). Functional interpretation of genetic variants using deep learning predicts impact of the epigenome. *bioRxiv*. <https://doi.org/10.1101/389056>.
- Huang, G., Liu, Z., Weinberger, K.Q., and van der Mateen, L. (2016). Densely Connected Convolutional Networks. *arxiv.org*, arXiv:1608.06993.
- Kelley, D.R., Snoek, J., and Rinn, J.L. (2016). Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999.
- Kim, J.H., and Kim, N. (2016). Signaling pathways in osteoclast differentiation. *Chonnam Med. J.* **52**, 12–17.
- Kuleshov, M.V., Jones, M.R., Rouillard, A.D., Fernandez, N.F., Duan, Q., Wang, Z., Koplev, S., Jenkins, S.L., Jagodnik, K.M., Lachmann, A., et al.

- (2016). Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* 44, W90–W97.
- Lanchantin, J., Singh, R., Lin, Z., and Qi, Y. (2016). Deep motif: visualizing genomic sequence classifications. arXiv:1605.01133v2 [cs.LG], 1–5. arXiv.org.
- van der Maaten, L., and Hinton, G. (2008). Visualizing data using t-SNE. *J. Mach. Learn. Res.* 9, 2579–2605.
- Maurano, M.T., Humbert, R., Rynes, E., Thurman, R.E., Haugen, E., Wang, H., Reynolds, A.P., Sandstrom, R., Qu, H., Brody, J., et al. (2012). Systematic localization of common disease-associated variation in regulatory DNA. *Science* 337, 1190–1195.
- Mella, M., Kauppila, J.H., Karihtala, P., Lehenkari, P., Jukkola-Vuorinen, A., Soini, Y., Auvinen, P., Vaarala, M.H., Ronkainen, H., Kauppila, S., et al. (2015). Tumor infiltrating CD8(+) T lymphocyte count is independent of tumor TLR9 status in treatment naïve triple negative breast cancer and renal cell carcinoma. *Oncoimmunology* 4, e1002726.
- Moskowitz, D.M., Zhang, D.W., Hu, B., Le Saux, S., Yanes, R.E., Ye, Z., Buenostro, J.D., Weyand, C.M., Greenleaf, W.J., and Goronzy, J.J. (2017). Epigenomics of human CD8 T cell differentiation and aging. *Sci. Immunol.* 2, eaag0192.
- Nagarsheth, N., Wicha, M.S., and Zou, W. (2017). Chemokines in the cancer microenvironment and their relevance in cancer immunotherapy. *Nat. Rev. Immunol.* 17, 559–572.
- Nair, S., Kim, D.S., Perricone, J., and Kundaje, A. (2019). Integrating regulatory DNA sequence and gene expression to predict genome-wide chromatin accessibility across cellular contexts. bioRxiv 605717, <https://doi.org/10.1101/605717>.
- Philip, M., Fairchild, L., Sun, L., Horste, E.L., Camara, S., Shakiba, M., Scott, A.C., Viale, A., Laeuer, P., Merghoub, T., et al. (2017). Chromatin states define tumour-specific T cell dysfunction and reprogramming. *Nature* 545, 452–456.
- Polak, P., Lawrence, M.S., Haugen, E., Stoletzki, N., Stojanov, P., Thurman, R.E., Garraway, L.A., Mirkin, S., Getz, G., Stamatoiyannopoulos, J.A., and Sunyaev, S.R. (2014). Reduced local mutation density in regulatory DNA of cancer genomes is linked to DNA repair. *Nat. Biotechnol.* 32, 71–75.
- Quang, D., and Xie, X. (2016). DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. *Nucleic Acids Res.* 44, e107.
- Rivas-Fuentes, S., Salgado-Aguayo, A., Pertuz Belloso, S., Gorocica Rosete, P., Alvarado-Vasquez, N., and Aquino-Jarquín, G. (2015). Role of chemokines in non-small cell lung cancer: angiogenesis and inflammation. *J. Cancer* 6, 938–952.
- Sarvaiya, P.J., Guo, D., Ulasov, I., Gabikian, P., and Lesniak, M.S. (2013). Chemokines in tumor progression and metastasis. *Oncotarget* 4, 2171–2185.
- Sheffield, N.C., Thurman, R.E., Song, L., Safi, A., Stamatoiyannopoulos, J.A., Lenhard, B., Crawford, G.E., and Furey, T.S. (2013). Patterns of regulatory activity across diverse human cell types predict tissue identity, transcription factor binding, and long-range interactions. *Genome Res.* 23, 777–788.
- Simonyan, K., and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv:1409.1556 [cs.CV]. ArXiv.org.
- Singh, S., Yang, Y., Poczos, B., and Ma, J. (2016). Predicting enhancer-promoter interaction from genomic sequence with deep neural networks. bioRxiv 85241, <https://doi.org/10.1101/085241>.
- Sudmant, P.H., Alexis, M.S., and Burge, C.B. (2015). Meta-analysis of RNA-seq expression data across species, tissues and studies. *Genome Biol.* 16, 287.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. 2016. Rethinking the inception architecture for computer vision. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2818–2826.
- TCGA Genome Characterization Center UNC (2017). Dataset: gene expression RNAseq - IlluminaHiSeq - TCGA.LUAD.sampleMap/HiSeqV2.
- TCGA (2018). The Cancer Genome Atlas Program. <https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>.
- Thorsson, V., Gibbs, D.L., Brown, S.D., Wolf, D., Bortone, D.S., Ou Yang, T.-H., Porta-Pardo, E., Gao, G.F., Plaisier, C.L., Eddy, J.A., et al. (2018). The immune landscape of cancer. *Immunity* 48, 812–830.e14.
- Toil RNA-seq Recompute, U. C. S. C. (2016). TCGA Pan-Cancer (PANCAN) - Gene Expression RNAseq - TOIL RSEM tpm. https://xenabrowser.net/datapages/?dataset=tcga_RSEM_gene_tpm&host=https://toil.xenahubs.net.
- Vaske, C.J., Benz, S.C., Sanborn, J.Z., Earl, D., Szeto, C., Zhu, J., Haussler, D., and Stuart, J.M. (2010). Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics* 26, i237–i245.
- Vivian, J., Rao, A.A., Nothhaft, F.A., Ketchum, C., Armstrong, J., Novak, A., Pfeil, J., Narkizian, J., Deran, A.D., Musselman-Brown, A., et al. (2017). Toil enables reproducible, open source, big biomedical data analyses. *Nat. Biotechnol.* 35, 314–316.
- Weinstock, M., and McDermott, D. (2015). Targeting PD-1/PD-L1 in the treatment of metastatic renal cell carcinoma. *Ther. Adv. Urol.* 7, 365–377.
- Xu, K., Ba, J., Kiro, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Yoshua, B. (2015). Show, attend and tell: neural image caption generation with visual attention. *Int. Conf. Mach. Learn.* 37, 2048–2057.
- Yasmin, R., Siraj, S., Hassan, A., Khan, A.R., Abbasi, R., and Ahmad, N. (2015). Epigenetic regulation of inflammatory cytokines and associated genes in human malignancies. *Mediators Inflamm.* 2015, 201703.
- Zhang, X., and Ho, S.M. (2011). Epigenetics meets endocrinology. *J. Mol. Endocrinol.* 46, R11–R32.
- Zhao, Y.T., Kwon, D.Y., Johnson, B.S., Fasolino, M., Lamonica, J.M., Kim, Y.J., Zhao, B.S., He, C., Vahedi, G., Kim, T.H., and Zhou, Z. (2018). Long genes linked to autism spectrum disorders harbor broad enhancer-like chromatin domains. *Genome Res.* 28, 933–942.
- Zhou, J., and Troyanskaya, O.G. (2015). Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods* 12, 931–934.

ISCI, Volume 20

Supplemental Information

Deep Learning Implicitly Handles Tissue

Specific Phenomena to Predict Tumor

DNA Accessibility and Immune Activity

Kamil Wnuk, Jeremi Sudol, Kevin B. Givechian, Patrick Soon-Shiong, Shahrooz Rabizadeh, Christopher Szeto, and Charles Vaske

Transparent Methods

Baseline tissue-specific dataset

All tissue-specific models described in this work were trained and evaluated following the exact procedure of the Basset network (Kelley et al., 2016), using DNase-seq peak data from 164 sample types obtained from ENCODE (Consortium, 2012) and Roadmap Epigenomics (Kundaje et al., 2015) projects.

Greedy merging of overlapping peaks across all DNase-seq data samples allowed us to create a universal set of potential accessibility sites. For each site, a binary vector was used to label its accessibility state in each of the 164 cell types. Data was then split by genomic site so that 70,000 peak locations were held out for validation, 71,886 for testing, and the remaining 1.8 million sites were used for training. The model input was a 600 base pair window of the DNA sequence centered at a site of interest, represented as one-hot encoding.

Baseline tissue-specific model implementation

We used TensorFlow (Abadi et al., 2015) to implement the Basset architecture for our baseline. We used Adam (Kingma and Ba, 2014) instead of RMSProp (Tieleman and Hinton, 2012) to optimize network parameters. We also found that use of a dynamic decay rate (that increased over the course of training) for updating moving averages in batch normalization (Ioffe and Szegedy, 2015) led to a model with competitive performance more quickly than when using a fixed decay. No other significant deviations from the original implementation were included.

When improving on the baseline model with convolutional layer factorizations, we focused experiments on factorizations that maintained the effective region of influence of the original layers and did not significantly increase the overall number of network parameters.

ENCODE DNase-seq and RNA-seq dataset

Data from the ENCODE project was initially collected at the start of 2017 for all cell or tissue types for which RNA-seq and DNase-seq measurements were both available. In order to capture a greater diversity, gene quantifications from RNA-seq files with the following ENCODE labels were collected: “RNA-seq”, “polyA mRNA”, “polyA depleted”, “single cell”. All files with “ERROR” audit flags were rejected. We kept files with “insufficient read depth,” and “insufficient read length” warnings. Despite being below ENCODE project standards, we believe the available read depths and lengths in warning situations were likely to be less of an issue when it comes to differentiating cell types (Conesa et al., 2016), and preferred to accept more potential noise in favor of a larger diversity of sample types.

The final step of data preparation involved assigning associations between specific RNA-seq and DNase-seq files within the same tissue type. In cases where there existed multiple exact matches of “biosample accession” identifiers between the two file types, associations were restricted to such exact matches. If exact accessions did not match, two file types were associated if it could be verified that they originated from the same tissue sample, cell line, or patient. This eliminated several tissue types for which no such correspondences existed. Both technical and biological replicates were treated as independent samples of the same tissue since we wanted to put the burden of learning non-invertible aspects of noise due to the measurement process on the neural network model.

The dataset was refined in late 2017, as several samples that had been part of our training and testing data were revoked by the ENCODE consortium due to quality concerns and updates. The final dataset consisted of 74 unique tissue types, distributed among partitions as discussed earlier (Table 2). The validation set was held constant, while the training and test sets included two variations.

We utilized the same greedy merge methodology described in Basset (Kelley et al., 2016) on all DNase-seq samples in our training sets to obtain a set of all potential sites of accessible DNA along the whole genome.

We used a fixed length of 600 base pairs (bp) centered at DHS peaks to define each site. Blacklisted sites at which measurements were suggested to be unreliable were excluded (Kundaje, 2016). This led to a total of 1.71 million sites of interest in the case of the held-out tissue data partition, and 1.75 million sites in the tissue overlap data partition. Using all sites across all available DNase-seq files, this produced a total 338.7 million training examples in the held-out tissue split.

As in other recent work on DNA-based prediction tasks (Alipanahi et al., 2015, Kelley et al., 2016, Singh et al., 2016, Quang and Xie, 2016) the sequence for each genomic site was obtained from human genome assembly hg19/GRCh37.

Training the expression informed model

During training data was balanced per batch due to a 14:1 ratio of negative to positive examples. Each batch sampled an equal amount of accessible and non-accessible sites without replacement, such that one pass through all available negative training examples constituted multiple randomly permuted passes through all positive training examples. In situations where a DNase-seq file had more than a single matching RNA-seq file, sites from that DNase-seq file were randomly assigned to one of the multitude of corresponding RNA-seq expression vectors each time they were selected for a training batch.

To generate a validation set that was manageable to evaluate frequently we selected 40,000 random samples from each of accessible and non-accessible sites per validation DNase-seq file. This resulted in a set of 440,000 validation examples that were used to estimate ROC AUC throughout training.

However, upon stopping we also evaluated prediction performance across whole genomes (all potential DHSs) of all validation samples (Supplemental Table S2). In cases where multiple RNA-seq file matches existed, predictions across the entire genome were evaluated once for every possible DNase-seq and RNA-seq file pair. Whole genome evaluation gave a better characterization of performance on the intended application, especially as captured by PR AUC, which is less misleading in the presence of data imbalance. Results on the test sets were evaluated across whole genomes following the same procedure.

The total number of examples (all sites across all samples) for validation was 20.5 million and 22.2 million for testing in the held-out tissue partition.

All RNA-seq expression data used to train and test models was in units of $\log_2(\text{TPM} + 1)$. There were many possible strategies for selecting the subset of genes for our input signature, but to initially avoid optimizing in this space, we relied on the prior work of the Library of Integrated Network-based Cellular Signatures (LINCS) and used their curated L1000 list of genes (LINCS, 2018). To ensure that the models could be applied later to cancer genomes in TCGA we converted all L1000 gene names into Ensembl gene identifiers and kept only those genes that were available in both ENCODE and TCGA TOIL RNA-seq files. After this refinement, our final input L1000 gene list consisted of 978 genes.

Expression informed model architecture and hyperparameters

We trained several alternative versions of our model and reported validation results over the course of training in Supplemental Figure S3 and Supplemental Figure S4.

The tissue-specific models demonstrated that multi-task outputs could share common convolutional layers and provide an accurate prediction of DNA accessibility across distinct sample types. Thus, we expected that if an input vector was discriminative of cell type it was likely to be sufficient to integrate it into the network after the convolutional layers. We evaluated adding a fully connected layer (depth = 500) before concatenating the vector of L1000 gene RNA-seq data to output from the convolutional layers, but found that it performed consistently worse (Supplemental Figure S3) than direct concatenation without the fully connected layer (Figure 1D).

Transfer learning consistently shortened the training time across model variants, and we found that using weights learned from the corresponding data partitions before final cleanup of revoked files was more effective on the validation set than was transfer of convolutional layer weights from the best tissue-specific model. However, our most impactful changes were increasing the batch size (from 128 to 512, and finally to 2048), and decreasing the learning rate (from 0.001 to 0.0001).

The tissue-specific models had multi-task outputs so that each training sample provided an information-rich gradient based on multiple labels for backpropagation. Since using RNA-seq inputs eliminated the need for multi-task outputs, each sample now only provided gradient feedback based on a single output. The batch size increase was intended to compensate for this change in output dimension to produce a more useful gradient for each batch.

The learning rate decrease, on the other hand, was guided by the observation that training was reaching a point of slow improvement before even a single full pass through all negative training examples. Our new dataset was also significantly larger than that used to train tissue-specific models.

We initialized our final expression-informed model (Figure 1D) with weights learned from the first iteration of the dataset, before erroneous revoked files were removed. In turn, those models were initialized with convolutional layer parameters from our best performing tissue-specific factorized convolutions model (Figure 1C). An effective batch size of 2048 was used for training (2 GPUs processing distinct batches of 1024), with an Adam (Kingma and Ba, 2014) learning rate of 0.0001 and a 0.25 fraction of positive to negative samples in every batch.

Expression informed model evaluation on ENCODE and genomic site annotations

ENCODE test set results were summarized in two ways: as a mean of AUC scores computed per whole genome sample (mean tissue type AUC in Tables 3 and 4), and as a single AUC score computed by considering predictions for all sites across all whole genome samples together (overall AUC in Tables 3 and 4). Only the latter was reported for performance analysis by genomic site type.

Two key sources were used to assign functional labels to accessibility prediction sites for performance breakdown. Exon, intragenic, and intergenic regions were derived from annotations defined by GENCODE v19 (Harrow et al., 2012). Promoter and promoter flank, and enhancer region annotations were obtained from the Ensembl Regulatory Build (Zerbino et al., 2015).

When investigating correlation of training similarity to test sample performance, since the modulating factor between predictions applied to different tissues is the input RNA-seq data, distance between test samples, t , and the training set, T , was computed as $d(t, T) = \min_{i \in T} \|\mathbf{r}_i - \mathbf{r}_t\|$, where \mathbf{r}_t is a test sample's vector of $\log_2(\text{TPM} + 1)$ expression levels for all L1000 genes.

Predicting DNA accessibility in TCGA

We applied our best expression informed model trained on the held-out tissue ENCODE partition to predict accessibility in TCGA. We restricted our predictions to promoter and promoter flank sites, since performance at those sites was high across all tests.

TOIL RNA-seq transcripts per million (TPM) gene expression data was used to obtain L1000 input gene signatures for all processed TCGA samples (Vivian et al., 2017, TOIL RNA-seq recompute, 2016). All expression values were converted from $\log_2(\text{TPM} + 0.001)$ to $\log_2(\text{TPM} + 1)$ before use.

For landscape views of accessibility and mutation impact analysis (Figure 3) we considered only samples with WGS available, and used mutation calls from an internal tool. For each sample site affected by at least one mutation, the change in predicted accessibility was computed before and after each mutation was applied, independently for SNPs and INDELS (Figure 3B). In order to apply t-SNE to generate the per-site landscape view (Figure 3D) we represented each site by a vector of binary accessibility decisions at that position across all selected TCGA samples with all mutations applied. All mutations were also applied when generating the per-patient t-SNE visualization (Figure 3F).

Accessibility in LUAD

To assess the uniqueness in perspective of accessibility versus RNA-seq, all LUAD samples for which we had WGS data were clustered into two groups via K-means. For this, four data sources were used: accessibility predictions for promoter and promoter flank sites, TOIL $\log_2(\text{TPM} + 1)$ RNA-seq gene expression data, HiSeqV2 $\log_2(\text{normalized count} + 1)$ RNA-seq gene expression data (TCGA Genome

Characterization Center, 2017) and TOIL $\log_2(\text{TPM} + 1)$ RNA-seq gene expression data for all genes in the L1000 gene set used as inputs to our expression informed model. For the first three datatypes, we clustered samples based on the top 5% most variable sites (for accessibility) or genes (for TOIL and HiSeqV2) across the LUAD cohort, following the logic that the most highly variable sites may highlight the most dramatic differentially active pathways. For the L1000 genes, clustering was based on the entire set of gene expression levels. To show the difference quantitatively between cluster assignments across data types we used adjusted mutual information (Vinh et al., 2010) (Figure 4B).

Exploration of pathway enrichment between the accessibility clusters was performed using Enrichr (Kuleshov et al., 2016). Genes for enrichment analysis were selected by first eliminating all genes below a standard deviation threshold of 0.33 (in TOIL data) across the LUAD cohort (in HiSeqV2 data the equivalently selected standard deviation threshold was 1.0). This threshold was selected to include the main peak of gene standard deviation and exclude the peak around zero (Supplemental Figure S5), comprised of genes with little change or very low levels of expression. All remaining genes were then compared with a two-sided t-test between the two clusters and p-values were adjusted with Benjamini Hochberg (BH) correction. Due to the low number of WGS samples in either cluster (21 samples in C0 and 20 samples in C1) a more permissive false discovery rate of 0.25 was chosen as the cutoff for differential expression. In TOIL data, this procedure returned 512 genes upregulated in C0 and 857 genes upregulated in C1. In HiSeqV2 data, the same process yielded 344 upregulated genes in C0 and 339 in C1.

For comparison of tumor mutational burden (TMB) across clusters, TMB was computed as the total count of missense and nonsense mutations in each WGS sample.

When the patient analysis set was expanded to include all LUAD samples without WGS mutation information, clustering based on promoter and promoter flank accessibility predictions was repeated with the same procedure as before (Figure 4D).

To investigate whether the accessibility space appeared continuous along the dimensions of most variance across LUAD we used Principal Component Analysis (PCA) applied to all promoter and flank accessibility predictions to project each sample onto the first three principal components (Figure 4F).

Correlating accessibility count with gene expression

Total accessibility count used to investigate gene correlations was computed as the total number of promoter and promoter flank sites predicted to be accessible after applying the binary decision threshold (at 80% precision) defined on ENCODE data. Again, only genes whose standard deviation was above 0.33 were considered for correlation analysis. Both Pearson and Spearman measures were evaluated, and the threshold for both measures was an absolute value above 0.4. All genes satisfying the threshold were analyzed for KEGG pathway enrichment with Enrichr (666 genes for Pearson correlation, and 418 genes for Spearman) (Supplemental Table S4 and S5)

Accessibility analysis in immune cell driven clusters

LUAD samples were clustered into two groups using K-means on vectors of lymphoid (21 cell types) and myeloid (13 cell types) xCell estimates (Aran et al., 2017), revealing a survival difference (Supplemental Figure S6C). We noticed that a plane orthogonal to the first principal component (PC1) partitioned cluster labels when xCell vectors were reduced to three dimensions with PCA (Supplemental Figure S6A and B). To exclude cases of near ambiguous label assignment and focus on more prominent differences we removed samples within a small margin at the midpoint between clusters (in PC1). Margin size was equal to the standard deviation of the smallest cluster in the PC1 dimension (Supplemental Figure S6D and E). After ignoring margin samples, the survival difference of patients between clusters increased in significance (logrank test $p = 6.7e-4$) (Supplemental Figure S6F).

Total methylation for all LUAD samples was computed as the sum of values at all sites measured by the Infinium HumanMethylation450 BeadChip, available from TCGA (TCGA, 2016). Total accessible site

count considered all promoter and promoter flank sites, with binary class assignment based on the 80% precision threshold (Figure 5B).

For further analysis of accessibility, only sites previously determined as facultative were considered and all with low standard deviation (< 0.135) across LUAD ($N = 512$) were eliminated, to ignore cohort specific constitutive sites with some tolerance for noise. The threshold was selected so that at minimum 10 accessibility values at a site had to be distinct from the site's values across the whole cohort (Supplemental Figure S7A).

Each accessibility prediction site was assigned to its nearest gene, according to distance in base pairs, as defined by GENCODE v19 (Harrow et al., 2012). We considered only accessibility sites within 50,000 base pairs as having a valid correspondence to a gene (Supplemental Figure S7B). Significantly differentiated accessibility sites were then used to vote for candidate upregulated genes in each cluster. A very conservative significance threshold (two sided t-test BH adj. $p < 1.0e-5$) was selected so as to only focus on the most striking accessibility differences. Each site was allowed to contribute a single vote to its corresponding gene according to the cluster in which the site was more accessible.

Genes with a consistent direction of upregulation votes were considered cluster-specific candidate genes to test for differential expression (532 genes in X0 and 2250 genes in X1). From the candidate genes for each cluster that also had significant (two sided t-test BH adj. $p < 0.01$) differential expression (190 in X0 and 835 in X1) we identified the group in each cluster that was consistent (123 in X0, and 536 in X1) and inconsistent (67 in X0, and 299 in X1) with the direction of increased accessibility. All four sets were then tested for KEGG pathway enrichment via Enrichr (Supplemental Table S6 and S7).

Predicting immune state from promoter and promoter flank accessibility

To train an ensemble of distinct models to discriminate immune hot from immune cold we used three fold cross validation; independently partitioning hot (X0 from immune cell based clustering) and cold (X1 from immune cell clustering) samples randomly to maintain an equal ratio across each fold. Training on

different random subsets of data enhanced robustness when dealing with training label uncertainty. Each classifier was an RBF kernel SVM with $C = 3.5$, and $\gamma = \frac{1}{N\sigma}$, where N was the number of features and σ was the standard deviation of feature values across the training set. Additionally, training samples were balanced by weights inversely proportional to class frequency, and Platt scaling (Platt, 1999) was used to obtain probability estimates from SVM classification. During ensemble classifier application we excluded all samples that did not have a mean probability of at least 0.5 for the ensemble's majority class prediction.

Input features to the classifier were binary accessibility predictions for a set of 484 sites comprised from the union of all immune hot (X0) sites consistent with gene expression and immune cold (X1) sites inconsistent with expression, as obtained from analysis of the xCell driven LUAD clusters. These sites were chosen both for their association with significant differences in expression of corresponding genes and the enrichment of those gene sets for immune relevant pathways.

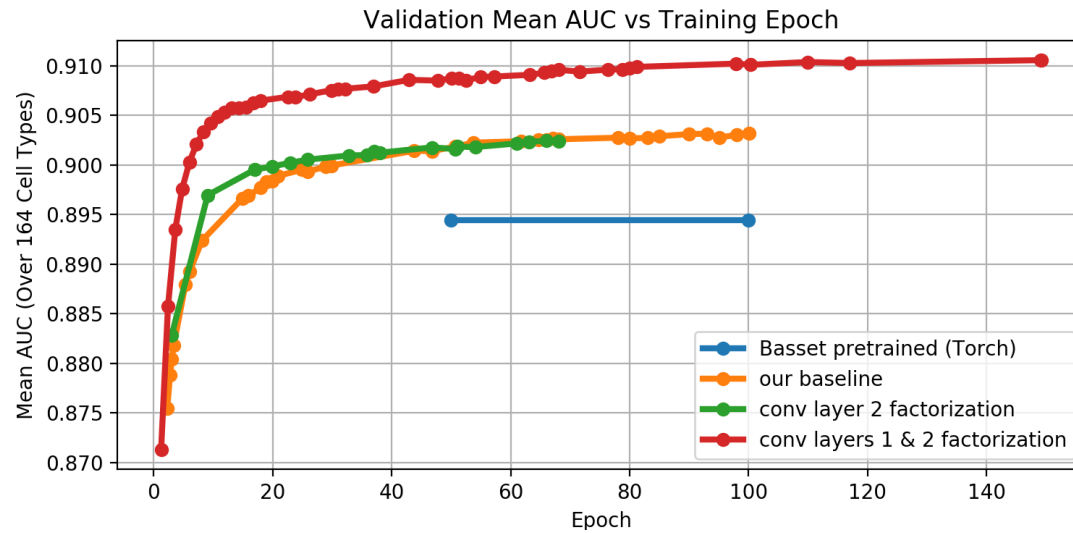
Expanding the application domain of the immune activity classifier to previously unprocessed TCGA cohorts involved first applying our expression informed convolutional neural network model to all promoter and promoter flank sites in the new data. As previously, when expanding our LUAD sample size to non-WGS data, we used only the reference genome (hg19/GRCh37) and TOIL $\log(\text{TPM} + 1)$ RNA-seq gene expression data for all predictions. Predictions that incorporated mutation information were included only for samples in our original six cohorts for which WGS was available.

Validating TCGA predictions with measured ATAC-seq peaks

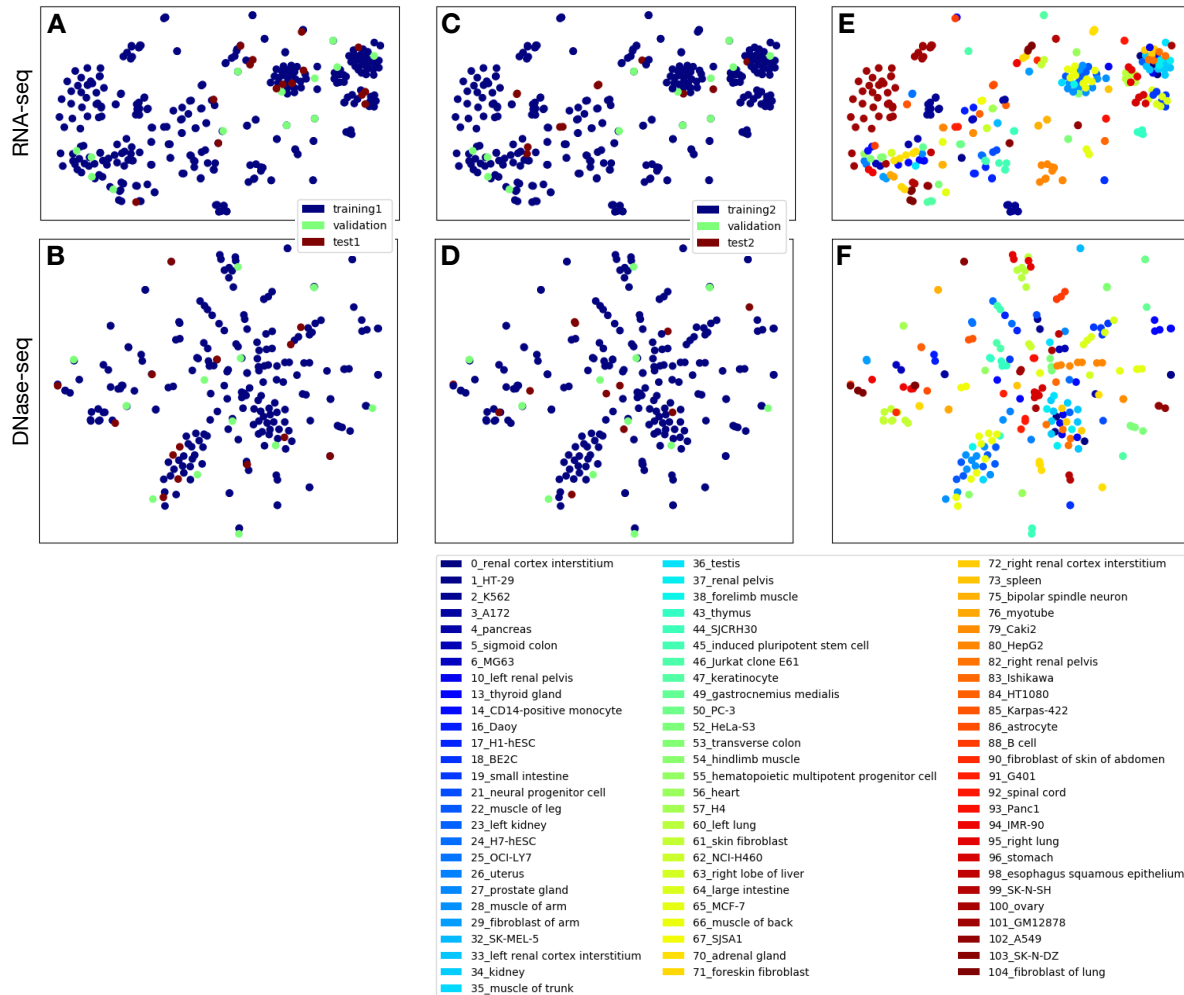
The list of pan-cancer peaks, TCGA sample identifiers, and the normalized ATAC-seq insertion counts within the pan-cancer peak set were obtained from supplemental material of the empirical investigation of chromatin accessibility in TCGA (Corces et al., 2018) at: <https://gdc.cancer.gov/about-data/publications/ATACseq-AWG> . To visualize distributions of ATAC-seq peaks for our clustering-based constitutive and facultative site labels, we computed the mean values within KIRC, KIRP, LUAD,

and LUSC cohorts individually for each pan-cancer peak that corresponded uniquely to our promoter and promoter flank sites with at least 70% overlap (Figure 7A). This was not restricted to TCGA samples for which we had also made predictions. But when looking at normalized count distributions for prediction decisions made by our accessibility classifier (Figure 7B,C) we did restrict analysis only to samples where both ATAC-seq was performed and our predictions were available. For every TCGA sample, each accessibility prediction site with a matching pan-cancer ATAC-seq peak contributed a single datum to the distribution plots. So every matched TCGA sample whose numbers are listed in Figure 7B,C contributed 61,342 data points. We validated matched samples for all cohorts to which we had previously applied our immune activity classifier with the exception of SARC, for which no ATAC-seq measurements were available, and GBM, for which none of the TCGA samples measured matched those for which we had run predictions. As in all previous TCGA analyses the classification decision threshold for binarizing accessibility predictions was based on an 80% precision (20% false discovery rate) threshold on the ENCODE held out tissue test set.

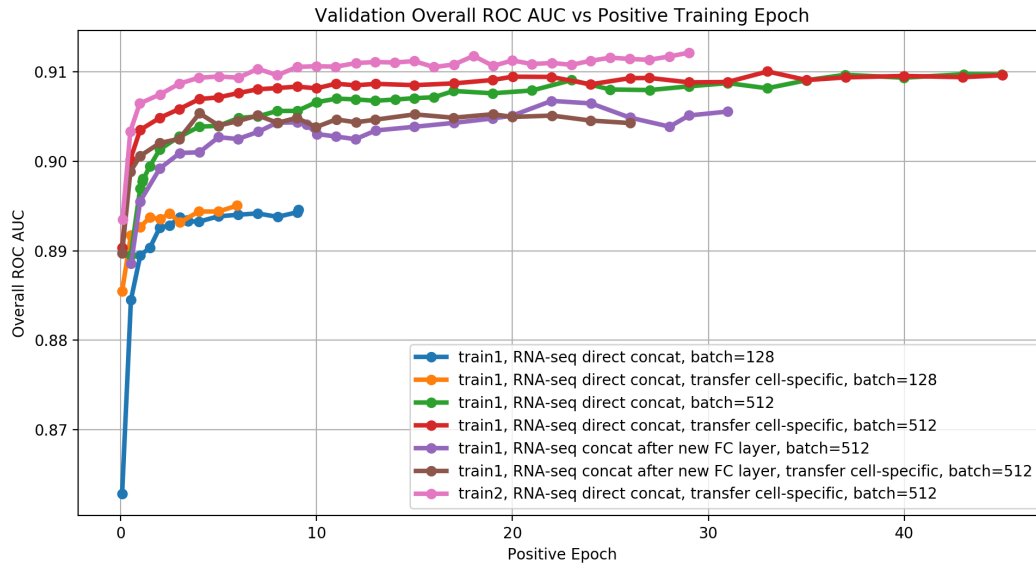
Supplemental Figures



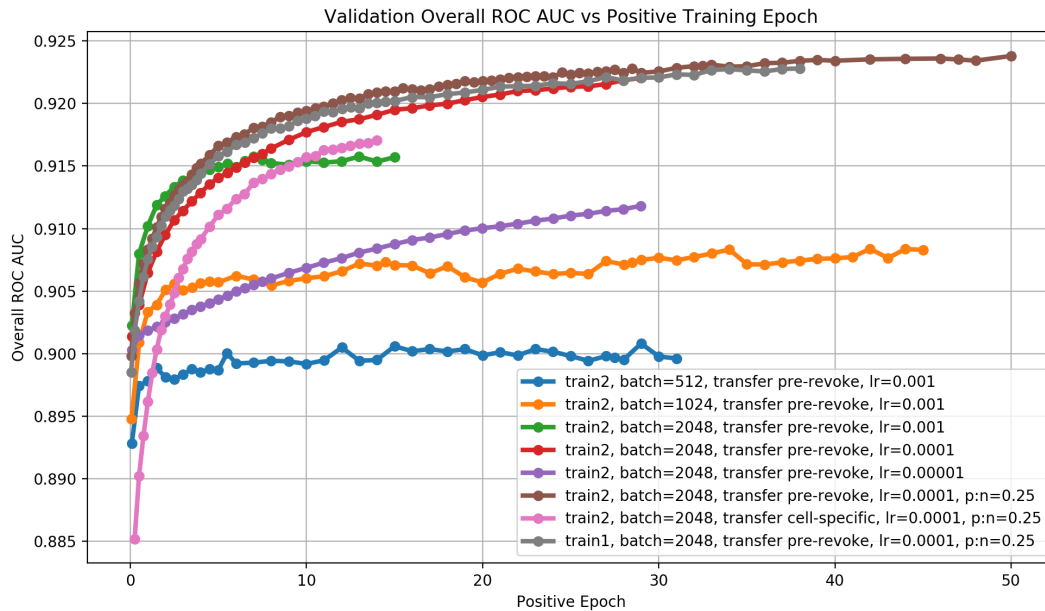
Supplemental Figure S1. Training of tissue-type-specific model architectures. Related to Figure 1 and Table 1. The mean ROC AUC across 164 cell types in the validation set versus training epoch is shown. The result obtained by the pre-trained model provided by the authors of Basset is shown for reference, but since the number of training epochs was not reported, an arbitrary range was selected for display. We explored independent factorization of the second convolutional layer of the baseline model, and achieved the best performance when both the first and second convolutional layers were factorized.



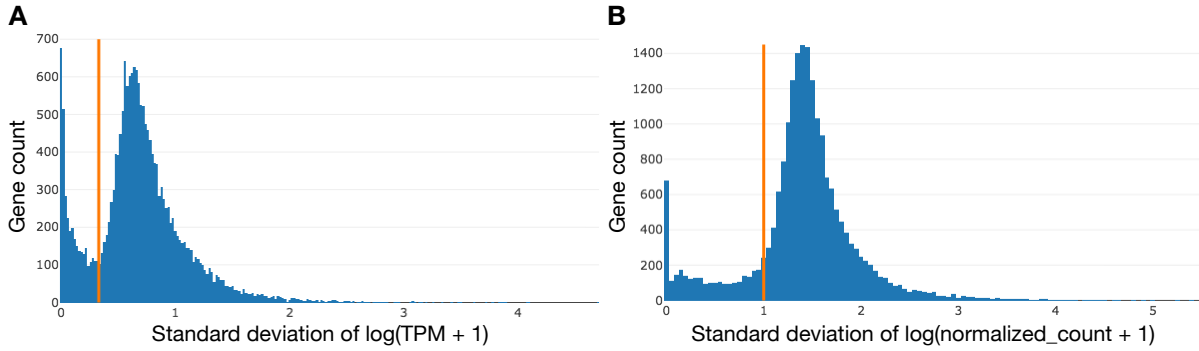
Supplemental Figure S2. t-SNE embedding of ENCODE dataset partitioning in RNA-seq and DNase-seq space. Related to Tables 2, 3, 4, and Figure 2. Sample distribution is illustrated by t-SNE embedding of the tissue overlap data partitions (training1 and test1) based on (A) RNA-seq $\log_2(\text{TPM} + 1)$ expression data and (B) DNase-seq peaks, as well as the held-out tissues data partitions (training2 and test2) based on (C) RNA-seq and (D) DNase-seq. The original ENCODE sample type labels are also shown for t-SNE embedded (E) RNA-seq and (F) DNase-seq samples, illustrating that samples of similar tissue or function often appear in proximity to each other across both data types.



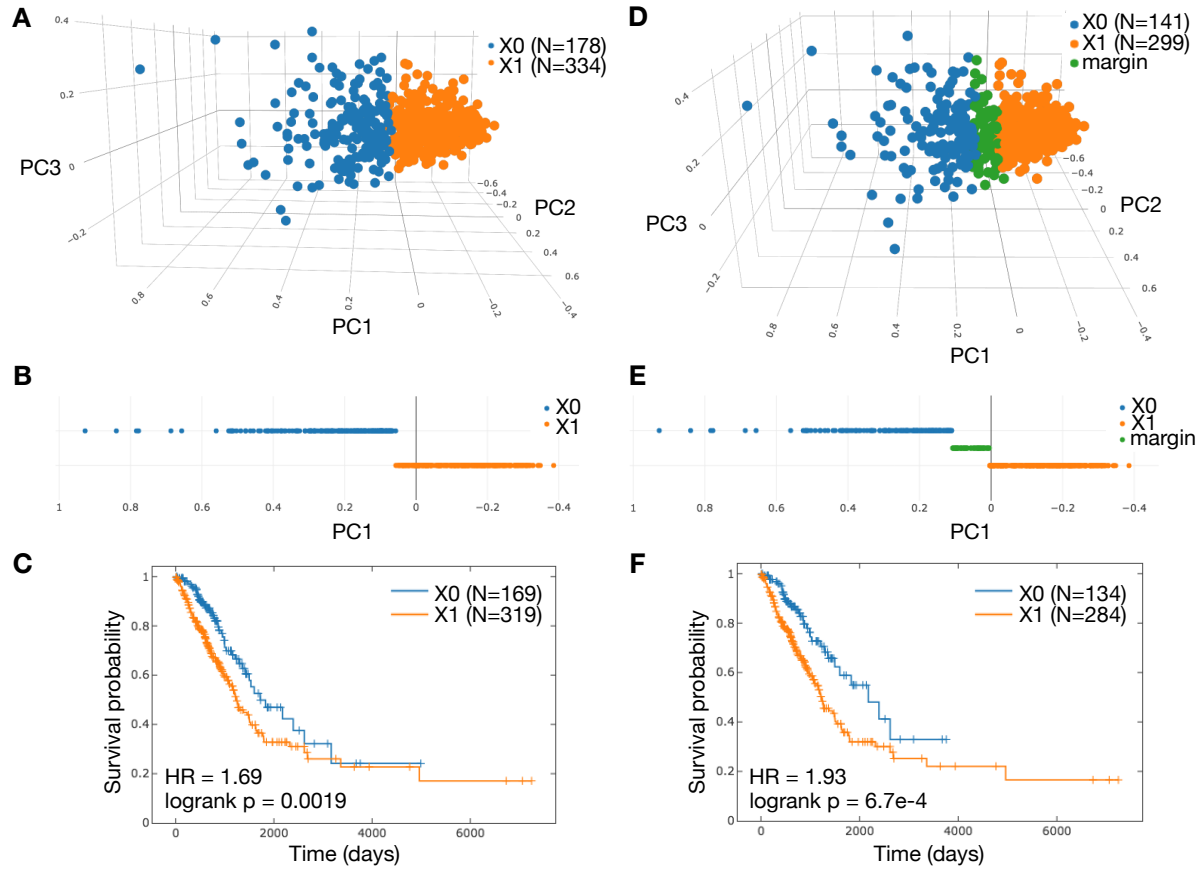
Supplemental Figure S3. Overall ROC AUC for the small validation set. Related to Figures 1, 2, and Tables 3, 4. The ROC AUC for the small validation set over number of passes through all positive examples (positive epochs) for several expression informed model architectures is shown. We experimented with adding a fully connected (FC) layer of depth 500 before concatenating (concat) gene expressions with outputs from the convolutional (conv) layers. However, increasing the batch size and initializing the convolutional layers with weights from our final tissue-specific model (transfer) improved performance most. Models trained on the tissue overlap set (train1) showed similar validation performance as those trained on the held-out tissue set (train2) with the same hyperparameters. This evaluation was done before the final dataset revision which revoked several suspected low quality samples, yet still provided valuable feedback for model selection.



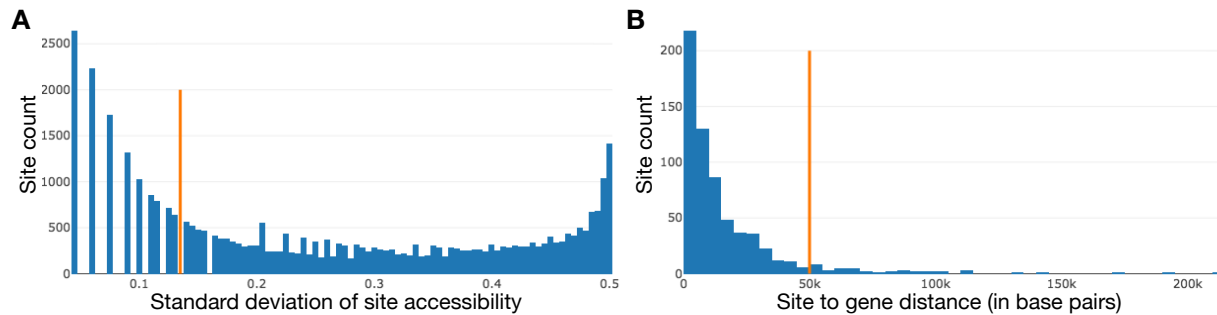
Supplemental Figure S4. Overall ROC AUC for the small validation set over positive training epochs for models trained after the final dataset revision. Related to Figures 1, 2, and Tables 3, 4. A further increase in batch size as well as a decreased learning rate (lr) led to additional significant improvements. Changing the fraction of positive samples per training batch (from p:n=0.5 to p:n=0.25) also slightly improved both ROC AUC as well as PR AUC in whole genome validation. Transfer of weights learned before final revoking of data (Figure S3) was a more effective initialization than weight transfer from our final tissue-specific model. Finally, we again confirmed that the same hyperparameters led to good validation performance across both training partitions: tissue overlap (train1) and held-out tissue (train2).



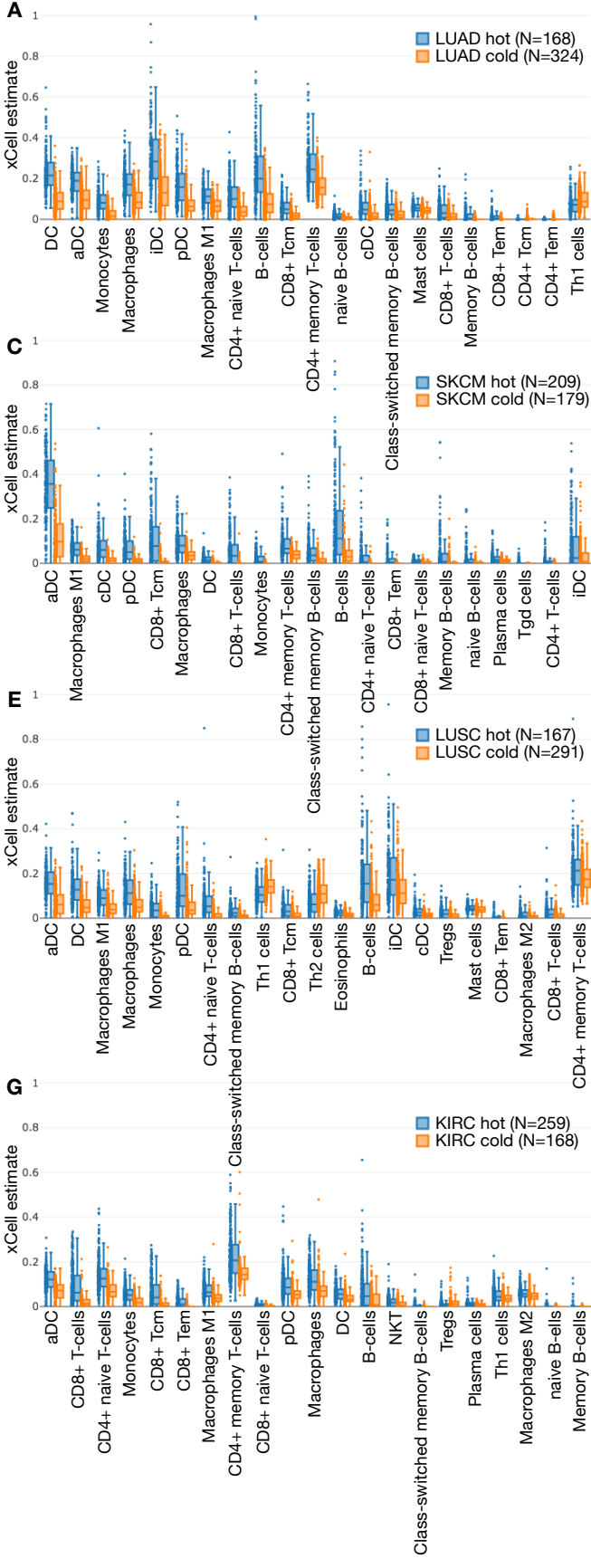
Supplemental Figure S5. Histograms of gene expression standard deviations. Related to Figure 4. (A) Histogram of gene expression standard deviations across LUAD WGS samples in TOIL RNA-seq $\log_2(\text{TPM} + 1)$ data along with the selected threshold (0.33) in orange is shown. (B) Gene expression standard deviations across the same samples in HiSeqV2 $\log_2(\text{normalized_count} + 1)$ data and the selected threshold (1.0) in orange (B) are also shown. In both cases the thresholds eliminate genes with little change across samples or very low levels of expression, and keep all genes that constitute the main peak of values.



Supplemental Figure S6. Adding a margin between immune cell based clusters in LUAD samples. Related to Figure 5. (A) All LUAD samples are shown with respect to the first three principal components (PC1-3) of their lymphoid and myeloid xCell estimates, colored according to labels assigned from k-means clustering. (B) Plotting the labeled data according to only the first principal component clearly shows the location of a separating plane between the clusters. The points excluded by introducing a margin at the scale of the smaller cluster's standard deviation along PC1 are shown (D) and (E). The impact on survival between X0 and X1 is shown by Kaplan-Meier plots before (C) and after (F) the margin was introduced. Kaplan-Meier plots are annotated with group size (N), logrank test p-values and hazard ratio (HR) based on a Cox proportional hazards (CoxPH) model regression using class assignment as the only explanatory variable.



Supplemental Figure S7. Histograms of accessibility and site to gene distance standard deviations. Related to Figure 5. (A) A histogram of the standard deviation of accessibility classifications in LUAD of promoter and promoter flank sites previously identified as facultative (40,823 sites) based on t-SNE across our initial set of six TCGA cohorts is shown. The threshold (< 0.135), in orange, identifies a subset of 25,093 sites facultative in LUAD. (B) The second histogram shows the site to nearest gene distances for all accessibility sites that also satisfy BH adjusted $p < 1.0e-5$ from a two-sided t-test between the LUAD immune cell driven clusters (3246 sites). Only sites within 50k base pairs (orange) were considered when voting for gene accessibility.



Supplemental Figure S8. Example immune cell and checkpoint gene distributions across predicted hot and cold tumors. Related to Figure 6. Examples of xCell estimates and checkpoint gene expression levels compared across multiple cohorts for tumors classified by our 3 SVM ensemble as hot or cold, with group size shown (N). All x-axis labels are ordered by significance based on a two-sided t-test between tumors classified as hot and cold. Only the top 21 most significant xCell estimates from lymphoid and myeloid cell categories are shown. (A) The LUAD xCell estimate and (B) checkpoint gene distributions demonstrate how application of the classifier affected significance ordering from the raw training data illustrated in Figure 5. (C,D) SKCM provides an example of a distinct cohort in which immune hot patients exhibited longer survival. (E,F) LUSC demonstrates one case where immune activity appears to have no effect. (G,H) Finally, KIRC shows a case where immune active patients have significantly worse survival, and also exhibits a curious lack of difference between levels of CD274 (also called PD-L1) and PDCD1LG2 (also called PD-L2) compared to other cohorts split by our accessibility based immune activity classifier.

Supplemental Tables

Supplemental Table S1. Number of DNase-seq files by tissue name per dataset partition, including tissue overlap (TO) and held-out tissue (HOT) sets. Related to Tables 2, 3, 4, and Figure 2.

ENCODE Tissue Name	Validation	Train (TO)	Test (TO)	Train (HOT)	Test (HOT)
renal cortex interstitium	0	3	0	3	0
HT-29	0	2	0	2	0
K562	0	2	0	2	0
A172	0	1	1	2	0
pancreas	0	1	0	1	0
MG63	0	2	0	2	0
left renal pelvis	0	3	1	4	0
thyroid gland	0	6	0	6	0
CD14-positive monocyte	1	1	0	1	0
Daoy	0	2	0	2	0
H1-hESC	0	2	0	2	0
BE2C	0	2	0	2	0
small intestine	0	3	2	5	0
neural progenitor cell	0	2	0	2	0
muscle of leg	0	7	0	7	0
left kidney	0	1	0	0	1
H7-hESC	0	3	0	3	0
OCI-LY7	0	2	0	0	2
uterus	0	1	0	1	0
prostate gland	0	1	0	0	1
muscle of arm	2	7	1	8	0
fibroblast of arm	1	1	0	1	0
SK-MEL-5	0	2	0	2	0
left renal cortex interstitium	0	4	0	4	0
kidney	0	4	0	4	0
muscle of trunk	0	1	0	1	0
testis	0	2	0	2	0
renal pelvis	0	3	0	3	0
forelimb muscle	0	0	1	1	0
thymus	0	6	0	6	0
SJCRH30	1	2	0	2	0
induced pluripotent stem cell	0	2	0	2	0
Jurkat clone E61	0	2	0	2	0
keratinocyte	0	1	1	2	0
PC-3	0	2	0	2	0
HeLa-S3	0	4	0	4	0
hindlimb muscle	0	1	0	0	1
hematopoietic multipotent progenitor cell	0	1	0	1	0
heart	0	2	0	2	0
H4	0	2	0	2	0
left lung	0	6	0	6	0
skin fibroblast	0	8	1	9	0
NCI-H460	0	2	0	2	0
large intestine	0	5	1	6	0
MCF-7	0	4	0	4	0
muscle of back	0	8	2	10	0
SJSA1	0	2	0	2	0
adrenal gland	0	5	1	6	0
foreskin fibroblast	1	1	0	1	0
right renal cortex interstitium	0	3	0	3	0
spleen	0	1	0	0	1
bipolar spindle neuron	0	2	0	2	0
myotube	0	2	0	2	0
Caki2	0	2	0	2	0
HepG2	0	4	0	4	0
right renal pelvis	1	3	0	3	0
Ishikawa	0	1	0	1	0
HT1080	0	2	0	2	0
Karpas-422	0	2	0	2	0
astrocyte	0	2	0	0	2
B cell	0	4	0	4	0
fibroblast of skin of abdomen	0	1	0	0	1
G401	0	2	0	0	2
spinal cord	1	1	0	1	0
Panc1	1	1	0	1	0
IMR-90	0	4	0	4	0
right lung	1	3	0	3	0
stomach	0	5	0	5	0
SK-N-SH	0	2	0	2	0
ovary	0	1	0	1	0
GM12878	0	2	0	2	0
A549	1	2	0	2	0
SK-N-IDZ	0	1	1	2	0
fibroblast of lung	0	5	1	6	0

Supplemental Table S2. Whole genome validation results for our expression-informed model trained on tissue overlap (TO) and held-out tissue (HOT) sets. Related to Tables 2, 3, and 4.

Sample tissue type	ROC AUC (TO)	PR AUC (TO)	ROC AUC (HOT)	PR AUC (HOT)
CD14-positive monocyte	0.888	0.559	0.889	0.563
muscle of arm	0.774, 0.959	0.654, 0.808	0.783, 0.960	0.671, 0.811
fibroblast of arm	0.898, 0.900	0.806, 0.809	0.898, 0.900	0.808, 0.811
SJCRH30	0.875	0.727	0.875	0.730
foreskin fibroblast	0.953	0.774	0.953	0.771
right renal pelvis	0.967	0.833	0.968	0.836
spinal cord	0.947	0.714	0.946	0.713
Panc1	0.957	0.713	0.958	0.711
right lung	0.958	0.781	0.958	0.782
A549	0.902	0.735	0.900	0.734
mean tissue type AUC	0.915	0.743	0.916	0.745
overall AUC	0.912	0.721	0.913	0.725

Supplemental Table S3. Enhancer results across held-out tissue test set whole genomes. Related to Table 6 and Figure 2.

Sample tissue type	ROC AUC	PR AUC
left kidney	0.934	0.737
OCI-LY7	0.845, 0.845, 0.850, 0.850	0.645, 0.643, 0.606, 0.605
prostate gland	0.817	0.490
hindlimb muscle	0.933	0.908
spleen	0.809	0.471
astrocyte	0.931, 0.898	0.967, 0.833
fibroblast of skin of abdomen	0.953	0.940
G401	0.640, 0.694	0.432, 0.270
overall AUC	0.870	0.732

Supplemental Table S4. Pathway enrichment (Enrichr) results with adjusted $p < 1.0e-4$ for all 418 genes correlated with total promoter and promoter flank accessibility in LUAD with $|\text{Spearman correlation}| > 0.4$. Related to Figure 4.

KEGG pathway	p	Adj. p	Z-score	Combined score
Osteoclast differentiation (hsa04380)	3.74e-17	7.45e-15	-1.86	70.51
TNF signaling pathway (hsa04668)	1.31e-10	1.30e-8	-1.89	42.96
Amoebiasis (hsa05146)	2.48e-9	1.64e-7	-1.82	36.03
Pathways in cancer (hsa05200)	2.33e-8	6.50e-7	-1.95	34.27
Tuberculosis (hsa05152)	6.68e-9	2.66e-7	-1.71	32.13
Pertussis (hsa05133)	4.72e-9	2.35e-7	-1.64	31.46
Regulation of actin cytoskeleton (hsa04810)	2.61e-8	6.50e-7	-1.75	30.52
NF-kappa B signaling pathway (hsa04064)	8.08e-9	2.68e-7	-1.62	30.28
Epstein-Barr virus infection (hsa05169)	1.28e-6	2.83e-5	-1.75	23.75
PI3K-Akt signaling pathway (hsa04151)	3.31e-6	5.21e-5	-1.78	22.48
Chemokine signaling pathway (hsa04062)	2.11e-6	4.01e-5	-1.69	22.15
Influenza A (hsa05164)	4.29e-6	6.10e-5	-1.63	20.18
Cytokine-cytokine receptor interaction (hsa04060)	3.40e-6	5.21e-5	-1.59	20.00
AGE-RAGE signaling pathway in diabetic complications (hsa04933)	8.58e-6	1.00e-4	-1.60	18.72
Jak-STAT signaling pathway (hsa04630)	6.09e-6	7.58e-5	-1.55	18.61
Staphylococcus aureus infection (hsa05150)	2.22e-6	4.01e-5	-1.43	18.56
Measles (hsa05162)	5.70e-6	7.56e-5	-1.50	18.17

Supplemental Table S5. Pathway enrichment (Enrichr) results with adjusted $p < 1.0e-6$ for all 666 genes correlated with total promoter and promoter flank accessibility in LUAD with $|\text{Pearson correlation}| > 0.4$. Related to Figure 4.

KEGG pathway	p	Adj. p	Z-score	Combined score
Osteoclast differentiation (hsa04380)	4.35e-22	9.66e-20	-1.86	91.69
Influenza A (hsa05164)	8.68e-13	9.63e-11	-1.94	53.97
TNF signaling pathway (hsa04668)	1.37e-12	1.01e-10	-1.86	50.83
Regulation of actin cytoskeleton (hsa04810)	5.54e-12	3.08e-10	-1.85	47.86
Hepatitis B (hsa05161)	9.68e-11	4.30e-9	-1.83	42.30
HTLV-I infection (hsa05166)	1.61e-10	5.95e-9	-1.82	41.00
Jak-STAT signaling pathway (hsa04630)	5.20e-10	1.44e-8	-1.75	37.47
Pathways in cancer (hsa05200)	1.76e-9	3.26e-8	-1.82	36.61
Chemokine signaling pathway (hsa04062)	7.17e-10	1.77e-8	-1.72	36.23
Epstein-Barr virus infection (hsa05169)	8.26e-10	1.83e-8	-1.73	36.09
Leukocyte transendothelial migration (hsa04670)	3.16e-10	1.00e-8	-1.57	34.40
Tuberculosis (hsa05152)	1.23e-9	2.48e-8	-1.56	32.02
Acute myeloid leukemia (hsa05221)	3.69e-9	6.29e-8	-1.55	30.04
Measles (hsa05162)	4.63e-9	7.15e-8	-1.53	29.37
Amoebiasis (hsa05146)	4.83e-9	7.15e-8	-1.48	28.35
Viral carcinogenesis (hsa05203)	2.30e-8	2.83e-7	-1.55	27.26
MAPK signaling pathway (hsa04010)	3.37e-8	3.94e-7	-1.54	26.49
B cell receptor signaling pathway (hsa04662)	1.42e-8	1.86e-7	-1.46	26.36
AGE-RAGE signaling pathway in diabetic complications (hsa04933)	3.67e-8	4.07e-7	-1.52	26.00
NF-kappa B signaling pathway (hsa04064)	1.01e-8	1.41e-7	-1.34	24.68
Focal adhesion (hsa04510)	7.25e-8	7.31e-7	-1.39	22.92
Fc gamma R-mediated phagocytosis (hsa04666)	6.71e-8	7.09e-7	-1.29	21.35

Supplemental Table S6. Top pathway enrichment (Enrichr) results for genes whose expression was consistent with increased accessibility in the immune active (X0) group of LUAD patients identified by xCell clustering. Related to Figure 5.

KEGG pathway	p	Adj. p	Z-score	Combined score
Focal adhesion (hsa04510)	0.000255	0.0355	-1.92	15.90
Osteoclast differentiation (hsa04380)	0.00135	0.0936	-1.84	12.15
PI3K-Akt signaling pathway (hsa04151)	0.00518	0.114	-1.99	10.47
Amoebiasis (hsa05146)	0.00339	0.114	-1.82	10.34
Acute myeloid leukemia (hsa05221)	0.00523	0.114	-1.88	9.86
Toxoplasmosis (hsa05145)	0.00610	0.114	-1.79	9.14
Proteoglycans in cancer (hsa05205)	0.00844	0.114	-1.81	8.62
Fc epsilon RI signaling pathway (hsa04664)	0.00851	0.114	-1.74	8.28
Renal cell carcinoma (hsa05211)	0.00784	0.114	-1.67	8.09
AMPK signaling pathway (hsa04152)	0.00725	0.114	-1.64	8.07
Rap1 signaling pathway (hsa04015)	0.00987	0.114	-1.68	7.78
Melanoma (hsa05218)	0.00957	0.114	-1.62	7.52

Supplemental Table S7. Pathway enrichment (Enrichr) results (adj. p < 0.05) for genes whose expression was inconsistent with increased accessibility in the immune cold (X1) group of LUAD patients identified by xCell clustering. Related to Figure 5.

KEGG pathway	p	Adj. p	Z-score	Combined score
Platelet activation (hsa04611)	2.24e-6	4.38e-4	-1.92	24.93
Inflammatory mediator regulation of TRP channels (hsa04750)	0.000112	0.0109	-1.99	18.07
Vascular smooth muscle contraction (hsa04270)	0.000449	0.0235	-1.82	14.00
Chemokine signaling pathway (hsa04062)	0.000529	0.0235	-1.85	13.96
cGMP-PKG signaling pathway (hsa04022)	0.000944	0.0235	-1.81	12.63
Focal adhesion (hsa04510)	0.000960	0.0235	-1.77	12.31
Intestinal immune network for IgA production (hsa04672)	0.000731	0.0235	-1.70	12.26
Cholinergic synapse (hsa04725)	0.00141	0.0247	-1.83	12.04
T cell receptor signaling pathway (hsa04660)	0.000960	0.0235	-1.69	11.72
Calcium signaling pathway (hsa04020)	0.00159	0.0247	-1.68	10.83
PI3K-Akt signaling pathway (hsa04151)	0.00197	0.0263	-1.73	10.78
Phospholipase D signaling pathway (hsa04072)	0.00148	0.0247	-1.64	10.70
Glutamatergic synapse (hsa04724)	0.00164	0.0247	-1.61	10.33
Pathways in cancer (hsa05200)	0.00275	0.0299	-1.66	9.77
ECM-receptor interaction (hsa04512)	0.00144	0.0247	-1.44	9.42
Long-term depression (hsa04730)	0.00202	0.0263	-1.42	8.79
Fc gamma R-mediated phagocytosis (hsa04666)	0.00273	0.0299	-1.41	8.30
Rap1 signaling pathway (hsa04015)	0.00461	0.0442	-1.51	8.12
Renin secretion (hsa04924)	0.00268	0.0299	-1.34	7.92
B cell receptor signaling pathway (hsa04662)	0.00474	0.0442	-1.36	7.27
Tuberculosis (hsa05152)	0.00544	0.0484	-1.29	6.74
Leishmaniasis (hsa05140)	0.00474	0.0442	-1.21	6.49