

Deep Learning Features Encode Interpretable Morphologies within Histological Images

Ali Foroughi pour

Jackson Laboratory

Brian White

Jackson Laboratory

Jonghanne park

Jackson Laboratory

Todd Sheridan

Jackson Laboratory

Jeffrey Chuang (✉ jeff.chuang@jax.org)

Jackson Laboratory

Research Article

Keywords:

Posted Date: February 14th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-865341/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

51 blurring, noise, and lossy image compression [4]. Tissue damage, image quality, and dataset-
52 specific artifacts have also been suggested to affect feature representation and prediction
53 accuracy of neural networks (see [1], [5], and [6]). Given the impact of such artifacts on deep
54 learning-based predictors, it is of critical importance to be able to decompose CNNs into
55 features that can be biologically interpreted.

56 The majority of models for visualizing, analyzing, and interpreting CNNs reveal “where”
57 a network is “looking” to make its prediction, rather than revealing “what” information in
58 the region of interest is important. Some methods output pixel patterns that affect the value
59 of a neuron in a deep network (see [7]). However, such techniques tend to output different
60 predictive regions, can be difficult to validate, or have been suggested to be “fragile”, i.e.
61 extremely sensitive to small perturbations of the image [8]. Optimizing conventional deep
62 learning techniques, such as self-attention, to identify regions informative of class labels
63 is a current theme in digital pathology (see [9, 10] for examples). While most methods
64 assess deep feature representations as a whole, recent work suggests deep learning features
65 cluster together and encode distinct morphologies [11]. Other recent works have focused
66 on visualizing individual deep learning features as heatmaps [12]. Finally, as the majority
67 of interpretation methods have focused on identifying regions predictive of class labels,
68 they requires a trained classifier and cannot be directly used in pipelines that employ
69 unsupervised feature learning. Unlike natural image analysis [13], biomedical image analysis
70 is complemented by additional data modalities, such as multiplexed imaging, single cell and
71 bulk sequencing, and clinical information citep ray2014information, kong2011integrative.
72 These data may aid in interpreting the deep feature representations of the H&E slide.
73 However, models integrating these diverse modalities are needed. The feasibility of doing so
74 is supported by work establishing the connection between modalities, for example by using
75 CNNs to predict expression values of specific genes from H&E images (see [14, 15, 16] for
76 examples). Because of the architectural complexity of CNNs, it has often been assumed that
77 CNN-based decompositions of images into features are not interpretable. However, there has
78 been little empirical study of this question, e.g. by testing whether CNN-derived features are
79 correlated with simple biological features such as gene expression values. In this work, we
80 investigate the interpretability of CNN-derived image features. Prior works ([17, 1]) have
81 referred to these by various names (e.g. features, fingerprints) whose use is not specific to
82 biological image analysis. For clarity and because they represent morphological features in
83 many ways analogous to genes, we refer to them as mones (i.e. “morphological genes”).
84 We find that mones share statistical similarities with gene expression data, and hence, a
85 mone can be conceptualized as an abstract gene with some expression value. Individual
86 mones have strong linear associations with phenotypic features, making them directly
87 interpretable, which we demonstrate in several analyses. We demonstrate that many mones
88 can distinguish cancer tumors from adjacent normal slides. These mones can be linearly
89 combined for reliable prediction of both pan-cancer and tissue-specific cancer phenotypes.
90 Mone-mone correlation analysis identifies clusters of highly correlated mones within cancer
91 types, and these correlations are strongly preserved among cancers from related tissues. The
92 similarity of mone values to gene expression data allows immediate use of many interpretable
93 bioinformatics tools and machine learning models to identify the underlying biology of
94 morphologies encoded by CNNs. For example, integrative mone-gene expression correlation
95 analysis reveals that collagen content and immune infiltration are linearly associated with
96 morphologies encoded by mones in several cancer types, and we confirm these relationships
97 by expert histopathological review. Our studies confirm individual deep learning features
98 encode distinct and identifiable morphology, and demonstrate the power of mones for
99 computationally deconstructing cancer images into interpretable biological features. The
100 linear analysis of individual deep learning features versus expression values or interpretable
101 morphology provides a simple and effective approach to interpreting deep learning models
102 in biological image analysis, notably without the need for a trained classifier.

103 **2. Results.** We analyzed the InceptionV3 [18] features of tiles from whole slide images
104 of The Cancer Genome Atlas (TCGA) [1] for 19 cancers (see Supplementary File 1 for the
105 full list). We used features derived from this architecture because predictive models based
106 on Inception have shown high accuracy for identifying phenotypes in prior studies [1]. We
107 hereafter denote each of the 2048 outputs of the global average pooling layer of the Inception
108 V3 network as mones (morphological genes). We use this terminology because mones have
109 analogies to genes with individual expression values. Tile level mones can be combined
110 to construct slide level mones (see methods). Unless otherwise stated, in the studies below
111 “mone” refers to a slide level characterization. Figure 1 provides an overview of interpretive
112 mone analyses and their connection to current interpretation techniques in the field.

113 *2.1. Individual mones differentiate phenotypes.* We first investigated to what extent
114 individual mones can differentiate phenotypes, focusing on TCGA tumor/normal slide
115 comparisons. We initially identified individual mones with significant differences in distribution
116 between breast cancer (BRCA) tumor and adjacent normal slides (>1800 mones were
117 statistically significant with $FDR < 5\%$). A clustermap of the top 100 such mones was
118 able to clearly separate these classes (see methods, Figure 2a, and Supplementary Figure
119 1, clustermap AUC=89%, rand score =96%, and adjusted rand score=85%). These results
120 were typical of mone behaviors in many tumor types – in any given tumor type, many mones
121 were able to separate frozen tumor from normal slides (see methods and Supplementary File
122 1). We applied several statistical methods to test the robustness of such mones. In each
123 cancer, at least 40% of mones were statistically significant irrespective of the statistical
124 test used ($FDR = 5\%$, see Supplementary File 1), and 75.7% of mones significant by at
125 least one method were significant by all tests (see methods, see Figure 2b). A smaller
126 subset of these mones showed strong effect sizes, as identified by optimal Bayesian filter
127 (OBF) [19] statistics (see Methods). 22% of all mone-cancer pairs met this criterion
128 based on distributional differences between tumor and normal slides (see Supplementary
129 Figure 2). As an example, mone 983 is a tumor marker with a distributional difference
130 between frozen tumor and adjacent normal breast cancer (BRCA) slides (Figure 2c). It
131 also behaves similarly in several other tumor types (Supplementary Figure 3). It strongly
132 correlates with cell density in frozen BRCA slides (see methods, see Figures 2d-e, Pearson
133 $r=0.69$, $p\text{-value} < 1e-200$, see Supplementary File 2) and significantly though with moderate
134 magnitude in FFPE BRCA tumor slides (Pearson $r=0.24$, $p\text{-value}=1.9e-13$). Mone 983 has
135 higher correlation with cell density in FFPE LUAD slides (Pearson $r=0.61$, $p\text{-value}=2.2e-$
136 34) than frozen LUAD slides (Pearson $r=0.28$, $p\text{-value}=4.1e-18$, see Supplementary File
137 2). To further test whether mones exhibit consistent behavior in different cancer types, we
138 analyzed the four cancer families of [1]: pan-GYN, pan-KIDNEY, pan-LUNG, and pan-
139 GI (see Supplementary File 3 for cancers in each family). More than half of the mones
140 with distributional differences in each cancer type also have distributional differences in
141 all cancers of the family (see Supplementary File 3). Although cancer-specific mones are
142 uncommon, such mones still show clear distributional differences between tumor and normal
143 (see Supplementary Figure 4). Interestingly, most mones distinguishing tumor from normal
144 in both LUAD and LUSC also distinguish between the LUAD and LUSC cancers (see
145 Supplementary File 3), suggesting quantitative values are important. Individual mones can
146 also distinguish frozen from FFPE slides (see Supplementary File 4). We observed the ability
147 of some mones to distinguish between tumor and normal is impacted by differences between
148 frozen and FFPE modalities (421 ± 112 across all cancers, see methods, see Supplementary
149 File 4), though the majority of mones behave similarly in frozen and FFPE.

150 *2.2. Mone clusters provide robust encodings of cancer phenotypes.* We next investigated
151 to what extent mones are independent or encode behaviors together. We did this by
152 calculating pairs of mones that were significantly correlated, for each cancer type. We
153 analyzed this first by restricting to frozen slides (to avoid Simpson’s paradox), and second
154 by combining frozen and FFPE slides (to avoid false correlations due to frozen-specific
155 artifacts). The results were generally robust between the two methods—a few cancers
156 had a non-trivial difference in the ratio of correlated mone pairs ascertained by the two
157 methods (ESCA, KICH, SARC, and PRAD $22.3\% \pm 2.22\%$), while the remaining cancers
158 had small differences ($7.8\% \pm 3.9\%$, see Supplementary File 6). Therefore we subsequently
159 analyzed frozen samples only unless otherwise specified. Overall, we found that correlated
160 mones are prevalent among tumor slides (see Figures 2f and 2g, and Supplementary File
161 6). For example, 83.9% and 88.7% of mone-mone pairs are correlated within LUAD and
162 within LUSC, respectively (see methods and Supplementary File 6). We observed similar
163 results for other cancers: $68.8\% \pm 13\%$ of all mone-pairs were statistically significant across
164 cancers (see Supplementary File 6). Remarkably, we observed that pairwise correlations are
165 preserved within cancer families – more than 45% of mone-pairs statistically significant in
166 one cancer are significant in all cancers of the family (pan-GYN, pan-KIDNEY, pan-LUNG,
167 and pan-GI, see Supplementary Figure 5). For example, the mone-mone correlations in lung
168 adenocarcinoma and lung squamous cell carcinoma are nearly identical (Figures 2f and 2g).
169 We also calculated mone-mone correlations in the normal slides associated with each cancer
170 type, hypothesizing that the difference in mone-mone correlations between tumor and normal
171 might be important to distinguishing tumor from normal images. However, these differential
172 mone correlations are weaker and less preserved across cancers (see Supplementary File 6
173 and Supplementary Figure 5).

174 Different cancers can be distinguished by different sets of mones. For example, while
175 mone 983 separates tumor and normal slides of BRCA, it does not differentiate COAD
176 tumor and normal slides (see Supplementary Figure 3). We identified 105 mones that are
177 highly correlated with mone 983 in BRCA (Null: $|r| \leq 0.5$, FDR= 0.1%, see methods),
178 among which 52, 14, and 29 also differentiate tumor from normal slides in COAD,
179 READ, and STAD, respectively (t-test FDR<0.1%). The first principal component of
180 these COAD-overlapping mones (explaining 41% of the variance across COAD samples)
181 strongly correlates with cell density in frozen slides (Cellpose cellularity estimates: Pearson
182 $r=0.22$, $p\text{-value}=2.2e-11$, HoverNet cellularity estimates: Pearson $r=0.43$, $P\text{-value}=1.9e-47$,
183 see Methods and Supplementary File 2). Thus the high cellularity in BRCA [20, 21] and
184 COAD [22, 23] involve incompletely overlapping mone sets.

185 *2.3. Linear models of mones can detect and distinguish tumors.* We investigated linear
186 models of mones for predicting phenotypes, as they allow direct interpretation of mone
187 values. We observed that linear models of mones can efficiently distinguish tumor from
188 adjacent normal slides, as well as the cancer type from which they are derived (19
189 cancers, 38 classes, see methods, see Figure 3a-c). We tested two linear models, multi-
190 class linear discriminant analysis [MLDA, One versus Rest (OVR)-AUC= $97.1\% \pm 4.6\%$,
191 see Supplementary File 7] and multinomial logistic regression with LASSO penalty (LR-
192 LASSO, OVR-AUC= $97.1\% \pm 4.2\%$, see Supplementary File 7). MLDA encodes mone
193 patterns indicative of class labels into a low dimensional space (i.e. the number of classes
194 - 1), yielding t-SNE visualizations with improved interpretability over naïve t-SNE (compare
195 Figures 3a,b and Supplementary Figures 6a,b). LR-LASSO, on the other hand, is a linear
196 model based on small mone sets, so its regression coefficients can be interpreted directly with
197 risks incurred by each mone. Although CNN methods typically use difficult-to-interpret fully
198 connected layers at the classification step, we found that efficiently designed linear models

199 can replace fully connected layers while still achieving high prediction AUCs. Combining
200 tumor probabilities of the LR-LASSO classifier we obtained a universal tumor detector with
201 extremely high AUC ($99.2\% \pm 0.12\%$, see methods), out-performing the fully deep learning
202 model of [24] (Reported AUC= 0.95 ± 0.02). Furthermore, LR-LASSO is effective at cross-
203 classification similar to CNNs with fully connected classification layers [1], i.e., LR-LASSO
204 trained to distinguish tumor/normal for one cancer type can distinguish tumor/normal for
205 other cancer types as well (Figure 3d and Supplementary Figure 6c). While the LR-LASSO
206 model has smaller average AUC (0.84) compared to the fully deep learning model of [1]
207 (0.88), logistic regression is more interpretable than a multi-layer perceptron. Our slide level
208 tumor detectors also produce meaningful tile level predictions. Independent review by our
209 pathology team supports most tumor regions having high tumor probability, and most non-
210 tumor regions have low tumor probability in these images, with the cases of misclassification
211 tending to be prediction of non-tumor regions to be tumor (see Supplementary Figure 7 for
212 examples). These results indicate that the LR-LASSO slide level tumor markers are effective
213 at the tile level.

214 *2.4. Mones have interpretable correlations with gene expression.* We next investigated
215 whether mones are linearly associated with gene expression values, as this could provide
216 transcriptional interpretability for mones (see Methods for data stratification and pre-
217 processing). We observed many mones significantly correlated with individual genes. Across
218 five cancers analyzed (OV, COAD, KIRC, LUAD, and LUSC) between 83 (LUSC) and 1797
219 (KIRC) mones were associated with at least one gene, whereas between 332 (LUSC) and
220 16474 (KIRC) genes were associated with at least one mone (methods, Supplementary File
221 8). We then analyzed several cases of particular interest.

222 *2.4.1. Mones encode collagen content.* We first used unsupervised analysis to study
223 clusters of correlated mone-gene pairs. We identified a cluster of highly correlated mones
224 and collagen genes in OV (Figure 4a). These mone values can be efficiently combined for
225 association with phenotype using PCA (PC-1 explains 63% of the variance). Histopathological
226 review by our pathology team confirmed tiles with high PC-1 values as typically rich in
227 collagen (see Figure 4d and Supplementary Figure 8), and tiles with low PC-1 values
228 as having low collagen but increased cellularity. These mone-gene associations may be
229 clinically relevant, as high expression of collagen genes correlates with multi-drug resistance
230 [25] and poor prognosis in ovarian cancers [26]. Mone 1062—one of the mones in the
231 identified cluster—was additionally highly correlated with ECM2, THBS1 and THBS2, which
232 have been suggested to play a significant role in ovarian cancer drug resistance and metastasis
233 [26, 27, 28].

234 *2.4.2. Mones encode immune infiltration.* Supervised correlation analysis using fixed
235 gene sets can also be used to test if mones encode morphological features associated with
236 a biological phenotype. We tested whether mones encode immune infiltration in pan-GI
237 (COAD, READ, and STAD) and pan-LUNG (LUAD and LUSC) cancers, with immune
238 related gene sets taken from [29] and [30] (see Methods and Supplementary File 9). We
239 identified 19 mones significantly correlating (some positive, some negative) with immune
240 related genes in pan-GI cancers (Figure 4b and Supplementary File 10). These mones
241 have significant correlations with each other ($|r|=0.56 \pm 0.13$), and 14 mones significantly
242 correlate with more than one immune gene. LAIR1, LCP2, MS4A4A, and CCR1 correlate
243 with 15, 9, 12, and 15 mones, respectively. All 19 mones are significantly correlated with
244 prior TCGA estimates of leukocyte fraction [31] ($FDR < 0.05, |r|=0.29 \pm 0.07$) and HoverNet
245 estimates of immune cell quantity from the H&E images (see methods, $FDR < 0.05, |r|=0.45 \pm$

0.12) in COAD. Mone 179 had strong positive correlations with both leukocyte fraction (246 $r=0.2$) and HoverNet estimates (247 $r=0.65$). Histopathological review confirmed that mone 179 differentiates amongst COAD tiles according to their level of immune infiltration (see 248 Figure 4e and Supplementary Figure 9). PC-1 of the 19 mones has significant correlation with 249 leukocyte fraction ($r=0.4$, $p\text{-value}=2.4e-35$) and HoverNet estimates ($r=0.17$, $p\text{-value}=2.5e-$ 250 17). Histopathological review validated that PC-1 also differentiates COAD tiles according to 251 differing levels of immune infiltration (see Supplementary Figure 10). Thus PC-1 efficiently 252 combines mones and provides a stronger separation than individual mones. We observed 253 similar correlations, but smaller in magnitude, between mones and local immune cytolytic 254 activity genes in lung cancers (Figure 3c and Supplementary File 10). We observed stronger 255 correlations in LUAD than LUSC (LUAD $|r|=0.20 \pm 0.05$, LUSC $|r|=0.12 \pm 0.05$). We 256 identified 31 mones significantly correlating (some positive, some negative) with immune 257 related genes in LUAD. 7 mones correlated with at least 3 genes (see Figure 4c). Only one 258 LUAD immune mone did not correlate with lymphocyte fraction ($FDR=0.05$). PC-1 of the 259 31 mones correlates with lymphocyte fraction ($r=0.28$, $p\text{-value}=2.2e-4$). Histopathological 260 review of LUAD slides based on PC-1 suggests LUAD tiles with high PC-1 show a strong 261 tumor infiltrating lymphocyte presence and have inflammation, while tiles with low PC-1 262 are typically not inflamed and show weaker immune infiltration (Figure 4f, Supplementary 263 Figure 11). 264

2.4.3. *Mones identify immunoglobulin gene expression in highly cellular colon* 265 *adenocarcinoma tumors.* Supervised correlation analysis can also clarify finer behaviors 266 within WSIs. For example, highly cellular COAD tumors typically show high expression 267 of immunoglobulin (IG) genes. We considered the 52 mones correlated with mone 983 and 268 which differentiate COAD tumor from normal (see section 2.2). The PC-1 dimension of these 269 mones significantly correlates with 87 genes in COAD ($FDR<0.05$, see Supplementary File 270 11 for the full gene list), as well with the average expression of these 87 genes ($r=0.33$, 271 $p\text{-value}=3.8e-12$). Immunoglobulin (IG) genes dominate this gene set, and we define their 272 average expression as a sample’s IG score. B-cells express immunoglobulin, and as expected 273 we observed a statistically significant correlation between the log normalized B-cell estimates 274 of [31] and IG score ($p\text{-value}=2.4e-12$). Interestingly, however, B cells comprise only a 275 small fraction of cells in each sample ($0.86\% \pm 1.5\%$), and we did not observe a significant 276 correlation between the B-cell estimates and mone PC-1 (correlation coefficient= -0.009 , $P\text{-}$ 277 $\text{value}=0.85$). Thus mone analysis suggests that IG expression is not due solely to B-cells. 278 This supports recent studies suggesting colon cancer cells themselves express IG genes (see 279 [32]). 280

2.5. *External Validation on CPTAC.* We externally validated mone patterns on CPTAC- 281 LUAD and CPTAC-LUSC. Mone 983 correlates with cell density in CPTAC-LUSC ($r=0.41$, 282 $p\text{-value}=1.1e-40$, Supplementary Figure 13) and CPTAC-LUAD ($r=0.21$, $p\text{-value}=1.8e-12$, 283 Supplementary Figure 13). Pathologist evaluation confirmed the ability of the PC-1 dimension 284 of TCGA-LUAD immune mones (see section 2.4.2) to separate CPTAC-LUAD slides by 285 immune activity. We also observed that HoverNet frequently mislabels lymphocytes as dead 286 cells (Supplementary Figure 13). Therefore, mone analysis has more robust behavior on 287 CPTAC-LUAD than HoverNet (Supplementary Figure 13). 288

3. Discussion. Mone analysis provides interpretability of deep neural networks through 289 the linear correlations of deep learning features against phenotypes and gene expression 290 profiles. We have empirically shown that mones efficiently encode strong morphological 291 features that can often be used to replace multi-layer perceptrons with robust and interpretable 292

293 linear classification models. Mone-analysis is flexible, can be used in a diverse set of
294 interpretation modalities, and can be applied to features engineered through various training
295 methodologies. Moreover, we have demonstrated that integrative mone-gene correlation
296 analysis can identify specific transcriptional processes from images, and verified these
297 through expert pathological review.

298 *3.1. Mones provide image-based interpretability.* Linear models based on mones have
299 several empirical and theoretical strengths for image analysis. Individual mones and small
300 mone clusters directly correlate with phenotypes (see section 2.1 and 2.3), enabling a simpler
301 interpretation of CNNs compared with methods that integrate all CNN features together. Most
302 interpretation models assume deep feature representations are complex and non-linear, and
303 therefore provide interpretability primarily through example regions identified by problem-
304 specific classifiers. On the other hand, our results demonstrate CNNs decompose H&E stained
305 images into interpretable features that linearly correlate with phenotypes, and the highly non-
306 linear feature representation assumption can be relaxed for interpreting CNNs trained on
307 H&E slides.

308 Some prior linear analyses of deep learning features have enabled partial interpretation
309 of CNNs trained on one cancer type (see [12] and [33]), but mones facilitate pan-cancer
310 interpretations. For example, our results demonstrate a pan-cancer mone can encode conserved
311 morphological features across multiple cancer types (see section 2.1), and a conserved morphological
312 feature can be encoded by distinct mones in different cancer types (see section 2.2). Moreover,
313 mones within strongly correlated clusters can be linearly combined to better identify shared
314 encoded morphology (see section 2.2 and [11]).

315 Mone analysis improves low dimensional visualization of large image sets, such as via
316 t-SNE plots. Pre-trained CNNs are universal feature extractors that encode morphological
317 features predictive of a multitude of labels (see [1], [34], and [35]). Because not all high-
318 dimensional complex relations can be easily embedded in 2D, mone analysis can be used as
319 a first phase of targeted search to identify the mones relevant to classes of interest (see section
320 2.3), allowing more accurate visualization of class separations in mone space.

321 Mone analysis guides classifier design by exploring statistical properties of the learned
322 features. For example, linear associations of correlated mone clusters with phenotypes suggests
323 utility of sparse linear models for reliable classification (see sections 2.3 and 2.5). Robustness
324 of linear models using a small number of mones (see sections 2.3 and 2.5) provides empirical
325 evidence for the theoretical results establishing robustness of sparse deep learning models
326 ([36],[37],[38]). Correlation analysis of mones with gene expression values is a powerful
327 approach for interpreting mones. We identified clusters of highly correlated mone-gene sets,
328 demonstrating clear connection of mones to the underlying genetics. Some recent studies
329 have used exhaustive sets of deep learning features to predict expression profiles (see [14, 15,
330 16]), but our work shows that small mone-gene clusters can be sufficient and provide simpler
331 interpretability. Both supervised and unsupervised analyses identify meaningful clusters (see
332 section 2.4). Supervised analysis using fixed gene sets is particularly interesting, as it enables
333 direct assessment of genes of interest via mones (see section 2.4.2).

334 *3.2. Interpretation without classifier training.* A notable advantage of the mone approach
335 is that it does not require a trained classifier, which is especially desirable when feature
336 engineering and classification are decoupled (e.g. transfer learning, unsupervised learning).
337 Interpretation models that require a trained classifier are restricted to the morphologies that
338 are predictive of a predetermined set of classes and are utilized by the classifier. Additionally,
339 different classification architectures applied to a fixed learned feature space may use different
340 features and morphologies for predicting a class label. Mone analysis does not require a

341 classifier to determine if the learned features encode a given phenotype (sections 2.4.2 and
342 2.4.3). It can be immediately applied to any learned feature space irrespective of training
343 methodology, and can be used in an unsupervised fashion to identify encoded morphologies
344 (section 2.4.1).

345 Optimization of training classifiers that are robust to stain differences across datasets is an
346 open research question. However, an advantage of mone analysis is that it can be more robust
347 to stain differences than pure classification models, as we observed for cell classification (see
348 section 2.5). It can also be used to analyze how stain differences affect feature representations
349 and a classifier’s ability to make reliable predictions.

350 *3.3. Mone analysis across architectures and data modalities.* Our analyses demonstrate
351 that mones provide an efficient and interpretable CNN embedding of image data, but a
352 caveat is that they have been restricted to Inception v3 mones. Architectures with fully
353 connected layers tend to increase non-linearity in feature representations. Therefore, models
354 that do not utilize multiple sequences of fully connected layers, such as Inception, are more
355 appropriate for linear mone analysis. For example, recent work suggests a small subset of
356 VGG19 features may also be interpreted via their direct association with phenotypes [12].
357 However, we believe Inception V3 mones are more appropriate for linear association studies
358 because they are the direct inputs to the classification layer. A few other studies have explored
359 correlations between deep autoencoder features and gene (see [33]) or protein expression [39]
360 for other architectures, but the relations between mones across architectures remains a broad
361 and open research topic. We have found that Inception V3 mone 983 in BRCA can be reliably
362 estimated via linear models using ResNet152V3 and DenseNet201 mones ($R^2 > 0.95$).
363 Furthermore, we have been able to convert Xception mones to Inception V3 mones using
364 autoencoders with reasonable accuracy ($R^2 \approx 0.5$). Recent work suggests InceptionV3 and
365 ResNet features are almost equivalent [40]. These studies suggest that integrative linear mone-
366 gene correlation analysis can be made effective across a range of deep learning architectures.

367 While this work focuses solely on H&E WSIs, we believe mone-based interpretation will
368 be valuable for extension to other spatial data types. Immunohistochemistry (IHC) images
369 are primed for rapid progress, as recent work has shown that several IHC markers can be
370 virtualized from H&E (see [41] and [42]). Generalizing mone analysis for other data types
371 such as spatial transcriptomics and multi-channel protein data is also an exciting and open
372 area, though new architectures will need to be explored to handle the high dimensionality of
373 such images. The interpretation of CNNs for these image types is a challenging but important
374 task, and we expect that integrative multi-modal multi-architecture mone analysis will be a
375 potent and informative approach.

376 **4. Methods.**

377 *4.1. Data acquisition and pre-processing.* The 20X H&E WSIs of TCGA are pre-
378 processed, tiled, and passed through an InceptionV3 model pre-trained on image-net data
379 in [1]. The cached 2048 global average pooling layer features of InceptionV3 (called mones
380 in this manuscript) were written to disk and analyzed. For each of these 2048 mones, the
381 median value across all tiles within a slide was computed to yield the slide-level mone value.

382 *4.2. Differential mone analysis.* Differential mone analysis identifies mones with
383 statistically significant distributional differences across classes. Welch’s t-test, Kolmogorov
384 Smirnov (KS) test, Wilcoxon Rank Sum (WRS) test, and optimal Bayesian Filter (OBF)
385 (see [19] for details) were used for statistical analysis. t-test, WRS test, and KS test use the
386 Benjamini-Hochberg procedure [43] for FDR correction. The scipy python package [44] was

387 used to implement t-test, KS test, and WRS test. The statsmodels [45] implementation of the
388 Benjamini-Hochberg procedure was used.

389 Minimal risk OBF (see [19] for details) identifies mones with posterior probabilities larger
390 than $1 - \alpha$, where α is the FDR rate. FDR-OBF (see [46] for details) outputs the feature set
391 that bounds the sample conditioned FDR by α . OBF can report the FDR of any arbitrary
392 feature set. Unless otherwise stated FDR-OBF is used. OBF uses Jeffrey's prior, assumes the
393 prior probability of a mone having distributional differences is 50% (to model no preference
394 on the identity of a mone, i.e., with or without distributional differences across classes), and
395 sets the normalization constant of the prior to 0.1. Mones with distributional differences
396 across classes are hereafter called markers, and mones without distributional differences
397 are called non-markers. The posterior probabilities of OBF can be used to estimate the
398 first two moments (mean and standard deviation) of the number of markers (see [46] for
399 details). Assuming mone identities (marker or non-marker) are independent across cancers,
400 the posterior probabilities can be multiplied to calculate the probability of a joint event.

401 OBF intrinsically computes the ratio between sample variance and weighted geometric
402 mean of class-conditioned variances, hereafter denoted by $a(m)$ for mone m . Similar to many
403 ANOVA-based analyses this ratio measures distributional differences and is closely related
404 to Bhattacharyya distance (see [47]). It converges to 1 for non-marker mones and converges
405 to larger values for markers. Larger $a(m)$ values denote larger distributional differences.
406 Assuming balanced samples of sizes 200 and 100 we compute the $a(m)$ values resulting in a
407 posterior of 0.95 as thresholds to distinguish moderate [$a(m)=1.088$] and strong [$a(m)=1.159$]
408 mones separating tumor from normal slides (see Supplementary Figure 3).

409 Structured multi-class OBF (see [48] for details) considers the four possible relations
410 (known as structures by OBF) between frozen normal, frozen tumor, and FFPE tumor
411 slides: (A) a mone does not differentiate between slides (prior probability=0.5), (B) a mone
412 has one distribution for frozen slides (both tumor and adjacent normal slides) and another
413 distribution for FFPE slides (prior probability=0.5/3), (C) a mone has one distribution for
414 tumor slides (both frozen and FFPE) and another distribution for frozen normal slides (prior
415 probability=0.5/3), and (D) a mone has one distribution for FFPE tumor slides and frozen
416 adjacent normal slides and another distribution for frozen tumor slides. Mones with structure
417 B for which frozen tumor and FFPE tumors lie on both sides of frozen adjacent normal slides
418 (based on mean values) are considered ineffective due to FFPE/frozen differences.

419 *4.3. Mone correlation analysis.* For each cancer type, we calculated correlations
420 between mones over all samples of the given cancer type. This analysis was done for each
421 cancer type. We use the Ledoit-Wolf shrinkage [49] implementation of the scikit-learn python
422 package [50] for computing covariance matrices. We then compute the correlation matrix
423 from the covariance matrix. We apply the Fisher transform to correlation coefficients and
424 approximate the null with its asymptotic Gaussian distribution. Benjamini-Hochberg [43]
425 procedure is used for FDR correction. We use seaborn package[51] with default values to
426 generate clustergrams. Correlation matrices are averaged to compute the pooled correlation
427 matrix. Statistically significant mone-mone correlations are referred to as "correlated mone
428 pairs."

429 Differential mone correlations denote the difference between the correlation coefficient of
430 tumor and normal slides, i.e., for each cancer, differential matrix is computed by subtracting
431 the correlation coefficient matrix of normal slides from the correlation coefficient matrix of
432 tumor slides. Differential mone correlation analysis uses the asymptotic Gaussian distribution
433 of the difference between Fisher transformed correlation coefficients to compute the p-values.
434 Statistically significant differential mone correlations are referred to as "differentially correlated
435 mone pairs."

436 4.4. *Linear classification models.* We implement MLDA and LR-LASSO using the scikit-
437 learn python package [50] with default values except we set $C = 100$ and use the “saga”
438 solver [52] in non-binary problems for LR-LASSO. We observed little sensitivity of AUCs
439 to the C values ranging from 1 to 1000, and hence use $C = 100$ throughout. We randomly
440 split data to train and test sets 10 times, which are then used to compute the mean and variance
441 of the AUCs. We use the Scipy package [44] to implement the Wilcoxon signed rank test.

442 4.5. *Cell segmentation and classification.* Cellpose [53] was used to segment and count
443 number of cells in BRCA, LUAD, and COAD tiles. Given the fixed magnification and tile size
444 the number of cells per tile captures tile level cell density. Median number of cells per tile was
445 used as slide level cell density index. HoverNet [54] was used to segment, count, and classify
446 nuclei within COAD tumor slides. HoverNet was executed using the pre-trained PanNuke
447 model ([55]), such that nuclei were classified into one of five types: neoplastic epithelial,
448 non-neoplastic epithelial, connective (including fibroblasts and endothelial), inflammatory
449 (including leukocytes, lymphocytes, and macrophages), and dead nuclei. Median number of
450 cells nuclei across tiles were used as cell density. Median number of predicted inflammatory
451 nuclei across tiles were used to characterize presence of immune cells.

452 4.5.1. *Integrative mone-gene analysis.* Gene expression data were downloaded from the
453 GDC portal [56]. We only used slide-gene expression pairs where both the slide and the
454 expression profile were from the same vial. Log normalized FPKMs were used. Genes with
455 zero counts in more than half the mone-gene pairs or expression standard deviation below
456 0.25 were removed. Given a set of mone-gene pairs, we stack the mone and gene vectors and
457 compute the covariance matrix using the Ledoit-Wolf shrinkage method [49] implemented in
458 the scikit-learn python package [50]. Correlation values are computed given the covariance
459 matrix similar to mone correlation analyses. Statistical significance tests are performed similar
460 to mone correlation analyses.

461 4.5.2. *Immune profiling and analysis.* Leukocyte fractions of TCGA samples were
462 obtained from [31]. All T-cell and B-cell categories were summed to obtain T-cell and B-
463 cell proportions, respectively. The fractions of T-cell and B-cells were summed to obtain
464 lymphocyte fractions. Log normalization of fractions were used throughout. Correlation
465 analysis of immune scores with mones and IG score was performed similar to mone
466 correlation analysis. B-cell percentages above 3% were removed for computing the B-cell
467 correlations as they were deemed outliers.

468 **Acknowledgements.** The authors would like to thank Jill Rubinstein for her feedback.
469 A. F. would like to thank The Jackson Laboratory for the JAX scholar award. J.H.C.
470 acknowledges support from grants R01CA230031 and P30CA034196.

471 **Conflict of interest.** The authors declare no conflict of interest.

REFERENCES

- 472 [1] Javad Noorbakhsh, Saman Farahmand, Ali Foroughi pour, Sandeep Namburi, Dennis Caruana, David Rimm,
473 Mohammad Soltanich-ha, Kourosh Zarringhalam, and Jeffrey H. Chuang.
474 Deep learning-based cross-classifications reveal conserved spatial behaviors within tumor histological
475 images.
476 *nature communications*.
477 [1](#), [2](#), [3](#), [5](#), [7](#), [8](#), [16](#)
- 478 [2] Riccardo Miotto, Fei Wang, Shuang Wang, Xiaoqian Jiang, and Joel T Dudley.
479 Deep learning for healthcare: review, opportunities and challenges.
480 *Briefings in bioinformatics*, 19(6):1236–1246, 2018.
481 [1](#)
- 482 [3] Shidan Wang, Donghan M Yang, Ruichen Rong, Xiaowei Zhan, Junya Fujimoto, Hongyu Liu, John Minna,
483 Ignacio Ivan Wistuba, Yang Xie, and Guanghua Xiao.
484 Artificial intelligence in lung cancer pathology image analysis.
485 *Cancers*, 11(11):1673, 2019.
486 [1](#)
- 487 [4] Samuel Dodge and Lina Karam.
488 Understanding how image quality affects deep neural networks.
489 In *2016 eighth international conference on quality of multimedia experience (QoMEX)*, pages 1–6. IEEE,
490 2016.
491 [2](#)
- 492 [5] Tejal Nair, Ali Foroughi pour, and Jeffrey H. Chuang.
493 The effect of blurring on lung cancer subtype classification accuracy of convolutional neural networks.
494 In *IEEE conference on bioinformatics and biomedicine*, pages 2987–2989. IEEE, 2020.
495 [2](#)
- 496 [6] Frederick M Howard, James Dolezal, Sara Kochanny, Jefree Schulte, Heather Chen, Lara Heij, Dezheng
497 Huo, Rita Nanda, Olufunmilayo I Olopade, Jakob N Kather, et al.
498 The impact of digital histopathology batch effect on deep learning model accuracy and bias.
499 *bioRxiv*, 2020.
500 [2](#)
- 501 [7] Maximilian Alber, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T Schütt, Grégoire
502 Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter-Jan Kindermans.
503 investigate neural networks!
504 *J. Mach. Learn. Res.*, 20(93):1–8, 2019.
505 [2](#), [16](#)
- 506 [8] Amirata Ghorbani, Abubakar Abid, and James Zou.
507 Interpretation of neural networks is fragile.
508 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688, 2019.
509 [2](#)
- 510 [9] Ming Y Lu, Melissa Zhao, Maha Shady, Jana Lipkova, Tiffany Y Chen, Drew FK Williamson, and Faisal
511 Mahmood.
512 Deep learning-based computational pathology predicts origins for cancers of unknown primary.
513 *arXiv preprint arXiv:2006.13932*, 2020.
514 [2](#)
- 515 [10] Miriam Hägele, Philipp Seegerer, Sebastian Lapuschkin, Michael Bockmayr, Wojciech Samek, Frederick
516 Klauschen, Klaus-Robert Müller, and Alexander Binder.
517 Resolving challenges in deep learning-based analyses of histopathological images using explanation
518 methods.
519 *Scientific reports*, 10(1):1–12, 2020.
520 [2](#)
- 521 [11] Ellery Wulczyn, David F Steiner, Melissa Moran, Markus Plass, Robert Reihs, Fraser Tan, Isabelle Flament-
522 Auvigne, Trissia Brown, Peter Regitnig, Po-Hsuan Cameron Chen, et al.
523 Interpretable survival prediction for colorectal cancer using deep learning.
524 *NPJ digital medicine*, 4(1):1–13, 2021.
525 [2](#), [7](#)
- 526 [12] Kevin Faust, Adil Roohi, Alberto J Leon, Emeline Leroux, Anglin Dent, Andrew J Evans, Trevor J Pugh,
527 Sangeetha N Kalimuthu, Ugljesa Djuric, and Phedias Diamandis.
528 Unsupervised resolution of histomorphologic heterogeneity in renal cell carcinoma using a brain tumor-
529 educated neural network.

- 530 *JCO Clinical Cancer Informatics*, 4:811–821, 2020.
 531 2, 7, 8
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei.
 533 Imagenet: A large-scale hierarchical image database.
 534 In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
 535 2
- [14] Benoît Schmauch, Alberto Romagnoni, Elodie Pronier, Charlie Saillard, Pascale Maillé, Julien Calderaro,
 536 Aurélie Kamoun, Meriem Sefta, Sylvain Toldo, Mikhail Zaslavskiy, et al.
 537 A deep learning model to predict rna-seq expression of tumours from whole slide images.
 538 *Nature communications*, 11(1):1–15, 2020.
 539 2, 7
 540
- [15] Alona Levy-Jurgenson, Xavier Tekpli, Vessela N Kristensen, and Zohar Yakhini.
 542 Spatial transcriptomics inferred from pathology whole-slide images links tumor heterogeneity to survival in
 543 breast and lung cancer.
 544 *Scientific reports*, 10(1):1–11, 2020.
 545 2, 7
- [16] Liviu Badea and Emil Stănescu.
 547 Identifying transcriptomic correlates of histology using deep learning.
 548 *PLOS ONE*, 15(11):1–30, 11 2020.
 549 2, 7
- [17] Rishi R Rawat, Itzel Ortega, Preeyam Roy, Fei Sha, Darryl Shibata, Daniel Ruderman, and David B Agus.
 551 Deep learned tissue “fingerprints” classify breast cancers by er/pr/her2 status from h&e images.
 552 *Scientific reports*, 10(1):1–13, 2020.
 553 2
- [18] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna.
 555 Rethinking the inception architecture for computer vision.
 556 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
 557 3
- [19] Ali Foroughi pour and Lori A. Dalton.
 559 Theory of optimal bayesian feature filtering.
 560 *Bayesian Analysis*, 15(4):1169–1197, 2020.
 561 3, 8, 9
- [20] Robert A Ambros and Ronald Trost.
 563 Cellularity in breast carcinoma.
 564 *American journal of clinical pathology*, 93(1):98–100, 1990.
 565 4
- [21] Kanji Tanaka, Daigo Yamamoto, Masanori Yamada, and Homa Okugawa.
 567 Influence of cellularity in human breast carcinoma.
 568 *The Breast*, 13(4):334–340, 2004.
 569 4
- [22] A Serrablo, P Paliogiannis, F Pulighe, S Saudi-Moro Moro, V Borrego-Estella, F Attene, F Scognamillo,
 571 and C Hörndler.
 572 Impact of novel histopathological factors on the outcomes of liver surgery for colorectal cancer metastases.
 573 *European Journal of Surgical Oncology (EJSO)*, 42(9):1268–1277, 2016.
 574 4
- [23] Frederikke Petrine Fliedner, Trine Bjørnbo Engel, Henrik H El-Ali, Anders Elias Hansen, and Andreas
 576 Kjaer.
 577 Diffusion weighted magnetic resonance imaging (dw-mri) as a non-invasive, tissue cellularity marker to
 578 monitor cancer treatment response.
 579 *BMC cancer*, 20(1):1–9, 2020.
 580 4
- [24] Jeonghyuk Park, Yul Ri Chung, Seo Taek Kong, Yeong Won Kim, Hyunho Park, Kyungdoc Kim, Dong-II
 582 Kim, and Kyu-Hwan Jung.
 583 Aggregation of cohorts for histopathological diagnosis with deep morphological analysis.
 584 *Scientific reports*, 11(1):1–11, 2021.
 585 5
- [25] Timon PH Buys, Raj Chari, Eric HL Lee, May Zhang, Calum MacAulay, Stephen Lam, Wan L Lam, and
 587 Victor Ling.
 588 Genetic changes in the evolution of multidrug resistance for cultured human ovarian cancer cells.

- 589 *Genes, Chromosomes and Cancer*, 46(12):1069–1079, 2007.
590 5
- 591 [26] Wei Zhang, Yi Liu, Na Sun, Dan Wang, Jerome Boyd-Kirkup, Xiaoyang Dou, and Jing-Dong Jackie Han.
592 Integrating genomic, epigenomic, and transcriptomic features reveals modular signatures underlying poor
593 prognosis in ovarian cancer.
594 *Cell reports*, 4(3):542–553, 2013.
595 5
- 596 [27] Su Li, Hua Li, Ying Xu, and Xiaomei Lv.
597 Identification of candidate biomarkers for epithelial ovarian cancer metastasis using microarray data.
598 *Oncology letters*, 14(4):3967–3974, 2017.
599 5
- 600 [28] Karolina Sterzyńska, Andrzej Klejewski, Karolina Wojtowicz, Monika Świerczewska, Marta Nowacka,
601 Dominika Kaźmierczak, Małgorzata Andrzejewska, Damian Rusek, Maciej Brązert, Jacek Brązert,
602 et al.
603 Mutual expression of *aldh1a1*, *lox*, and collagens in ovarian cancer cell lines as combined cscs-and ecm-
604 related models of drug resistance development.
605 *International journal of molecular sciences*, 20(1):54, 2019.
606 5
- 607 [29] Michael S Rooney, Sachet A Shukla, Catherine J Wu, Gad Getz, and Nir Hacohen.
608 Molecular and genetic properties of tumors associated with local immune cytolytic activity.
609 *Cell*, 160(1-2):48–61, 2015.
610 5
- 611 [30] Luigi Cari, Francesca De Rosa, Maria Grazia Petrillo, Graziella Migliorati, Giuseppe Nocentini, and Carlo
612 Riccardi.
613 Identification of 15 t cell restricted genes evaluates t cell infiltration of human healthy tissues and cancers
614 and shows prognostic and predictive potential.
615 *International journal of molecular sciences*, 20(20):5242, 2019.
616 5
- 617 [31] Vésteinn Thorsson, David L Gibbs, Scott D Brown, Denise Wolf, Dante S Bortone, Tai-Hsien Ou Yang,
618 Eduard Porta-Pardo, Galen F Gao, Christopher L Plaisier, James A Eddy, et al.
619 The immune landscape of cancer.
620 *Immunity*, 48(4):812–830, 2018.
621 5, 6, 10
- 622 [32] Zi-Han Geng, Chun-Xiang Ye, Yan Huang, Hong-Peng Jiang, Ying-Jiang Ye, Shan Wang, Yuan Zhou, Zhan-
623 Long Shen, and Xiao-Yan Qiu.
624 Human colorectal cancer cells frequently express *igg* and display unique *ig* repertoire.
625 *World journal of gastrointestinal oncology*, 11(3):195, 2019.
626 6
- 627 [33] Jordan T Ash, Gregory Darnell, Daniel Munro, and Barbara E Engelhardt.
628 Joint analysis of expression levels and histological images identifies genes associated with tissue
629 morphology.
630 *Nature communications*, 12(1):1–12, 2021.
631 7, 8
- 632 [34] Yu Fu, Alexander W Jung, Ramon Viñas Torne, Santiago Gonzalez, Harald Vöhringer, Artem Shmatko,
633 Lucy R Yates, Mercedes Jimenez-Linan, Luiza Moore, and Moritz Gerstung.
634 Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis.
635 *Nature Cancer*, 1(8):800–810, 2020.
636 7
- 637 [35] Jakob Nikolas Kather, Lara R Heij, Heike I Grabsch, Chiara Loeffler, Amelie Echle, Hannah Sophie Muti,
638 Jeremias Krause, Jan M Niehues, Kai AJ Sommer, Peter Bankhead, et al.
639 Pan-cancer image-based detection of clinically actionable genetic alterations.
640 *Nature Cancer*, 1(8):789–799, 2020.
641 7
- 642 [36] Subutai Ahmad and Luiz Scheinkman.
643 How can we be so dense? the robustness of highly sparse representations.
644 *CoRR*, vol. abs, page 11, 2019.
645 7
- 646 [37] Adnan Siraj Rakin, Zhezhi He, Li Yang, Yanzhi Wang, Liqiang Wang, and Deliang Fan.
647 Robust sparse regularization: Simultaneously optimizing neural network robustness and compactness.

- 648 *arXiv preprint arXiv:1905.13074*, 2019.
649 7
- 650 [38] Yan Sun, Wenjun Xiong, and Faming Liang.
651 Sparse deep learning: A new framework immune to local traps and miscalibration.
652 *Advances in Neural Information Processing Systems*, 34, 2021.
653 7
- 654 [39] Luke Ternes, Mark Dane, Marilynne Labrie, Gordon Mills, Joe Gray, Laura Heiser, and Young Hwan Chang.
655 Me-vae: Multi-encoder variational autoencoder for controlling multiple transformational features in single
656 cell image analysis.
657 *bioRxiv*, 2021.
658 8
- 659 [40] David McNeely-White, J Ross Beveridge, and Bruce A Draper.
660 Inception and resnet features are (almost) equivalent.
661 *Cognitive Systems Research*, 59:312–318, 2020.
662 8
- 663 [41] Christopher R Jackson, Aravindhan Sriharan, and Louis J Vaickus.
664 A machine learning algorithm for simulating immunohistochemistry: development of sox10 virtual ihc and
665 evaluation on primarily melanocytic neoplasms.
666 *Modern Pathology*, pages 1–11, 2020.
667 8
- 668 [42] Zhaoyang Xu, Carlos Fernández Moro, Béla Bozóky, and Qianni Zhang.
669 Gan-based virtual re-staining: a promising solution for whole slide image analysis.
670 *arXiv preprint arXiv:1901.04059*, 2019.
671 8
- 672 [43] Yoav Benjamini and Yosef Hochberg.
673 Controlling the false discovery rate: a practical and powerful approach to multiple testing.
674 *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
675 8, 9
- 676 [44] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni
677 Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stefan J. van der Walt, Matthew Brett,
678 Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern,
679 Eric Larson, CJ Carey, Ihan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef
680 Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R Harris, Anne M. Archibald,
681 Antonio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors.
682 Scipy 1.0: Fundamental algorithms for scientific computing in python.
683 *Nature Methods*, 2020.
684 8, 10
- 685 [45] Skipper Seabold and Josef Perktold.
686 statsmodels: Econometric and statistical modeling with python.
687 In *9th Python in Science Conference*, 2010.
688 9
- 689 [46] A. Foroughi Pour and L. A. Dalton.
690 Bayesian error analysis for feature selection in biomarker discovery.
691 *IEEE Access*, 7:127544–127563, 2019.
692 9
- 693 [47] Ali Foroughi pour and Lori A Dalton.
694 Optimal bayesian filtering for biomarker discovery: Performance and robustness.
695 *IEEE/ACM transactions on computational biology and bioinformatics*, 17(1):250–263, 2018.
696 9
- 697 [48] Ali Foroughi pour and Lori A Dalton.
698 Biomarker discovery via optimal bayesian feature filtering for structured multiclass data.
699 In *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and*
700 *Health Informatics*, pages 331–340, 2018.
701 9
- 702 [49] Olivier Ledoit and Michael Wolf.
703 A well-conditioned estimator for large-dimensional covariance matrices.
704 *Journal of multivariate analysis*, 88(2):365–411, 2004.
705 9, 10

- 706 [50] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,
707 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
708 E. Duchesnay.
709 Scikit-learn: Machine learning in Python.
710 *Journal of Machine Learning Research*, 12:2825–2830, 2011.
711 [9](#), [10](#)
- 712 [51] Michael Waskom and the seaborn development team.
713 mwaskom/seaborn, September 2020.
714 [9](#)
- 715 [52] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien.
716 Saga: A fast incremental gradient method with support for non-strongly convex composite objectives.
717 *arXiv preprint arXiv:1407.0202*, 2014.
718 [10](#)
- 719 [53] Carsten Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu.
720 Cellpose: a generalist algorithm for cellular segmentation.
721 *Nature Methods*, 18(1):100–106, 2021.
722 [10](#), [17](#)
- 723 [54] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and
724 Nasir Rajpoot.
725 Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images.
726 *Medical Image Analysis*, 58:101563, 2019.
727 [10](#)
- 728 [55] Jevgenij Gamper, Navid Alemi Koohbanani, Ksenija Benes, Simon Graham, Mostafa Jahanifar, Syed Ali
729 Khurram, Ayesha Azam, Katherine Hewitt, and Nasir Rajpoot.
730 Pannuke dataset extension, insights and baselines.
731 *arXiv preprint arXiv:2003.10778*, 2020.
732 [10](#)
- 733 [56] Robert L Grossman, Allison P Heath, Vincent Ferretti, Harold E Varmus, Douglas R Lowy, Warren A Kibbe,
734 and Louis M Staudt.
735 Toward a shared vision for cancer genomic data.
736 *New England Journal of Medicine*, 375(12):1109–1112, 2016.
737 [10](#)
- 738 [57] Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai
739 Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu.
740 Attention mechanisms in computer vision: A survey.
741 *arXiv preprint arXiv:2111.07624*, 2021.
742 [16](#)

743 **Figures.**

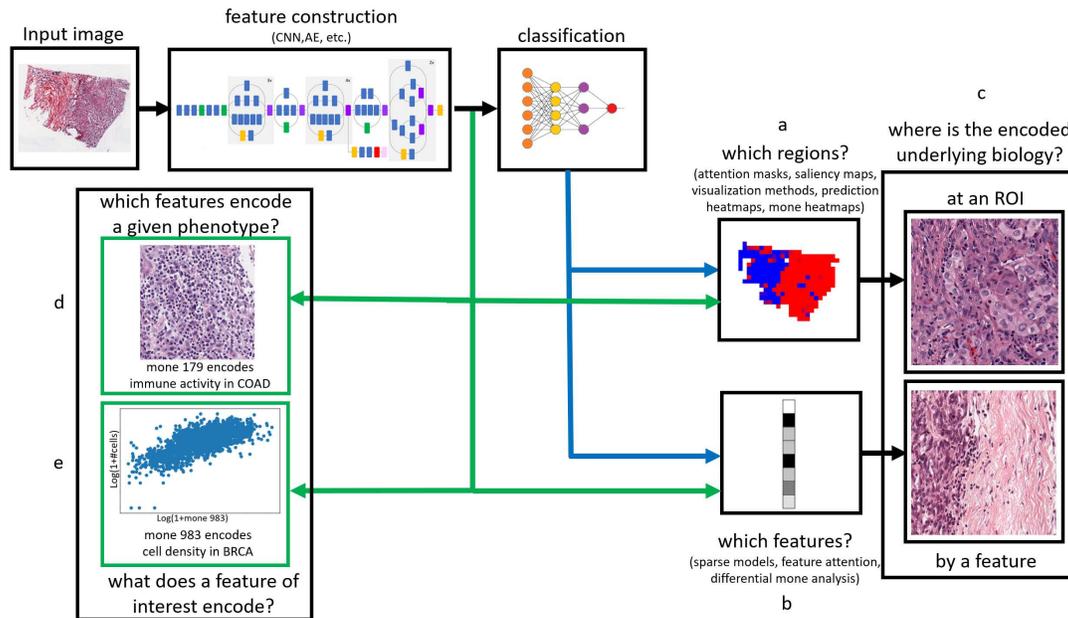


Figure 1: An overview of interpretation methods in deep learning. Blue arrows denote methods that require a trained classifier, and green arrows denote methods that do not require a trained classifier. (a) Several methods identify regions which drive the network’s prediction. These masks can be generated by the network, e.g. spatial self-attention [57], or as a post-process via visualization methods [7] such as GradCam, or prediction heatmaps [1]. Heatmaps of individual mones and mone-based classifiers can be used to detect predictive regions. (b) Channel attention [57] and sparse models, including sparse mone-based classifiers, identify subsets of features that are predictive of class labels. Differential mone analysis identifies discriminative features without training of a classifier. (c) Methods in (a) and (b) are typically used to select example image regions that have high attention, affect predictions the most, or affect the value of a feature. Mone analysis can be used to (d) identify features that encode a given phenotype of interest and (e) identify the morphology a feature of interest encodes without training a classification model.

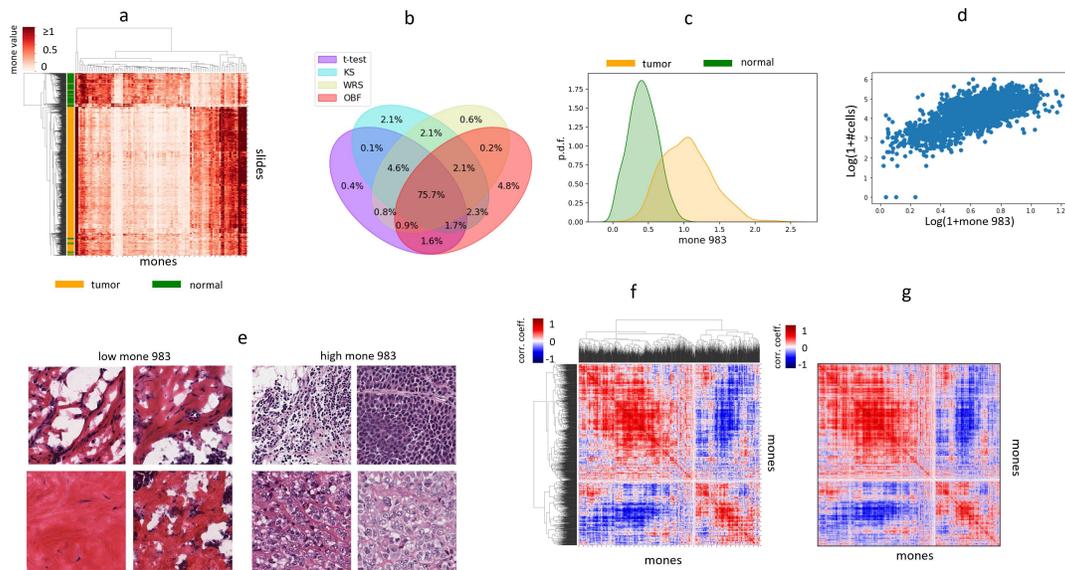


Figure 2: Individual mones and mone pairs encode and distinguish phenotypes. (a) Clustermap of BRCA slides using the top100 mones differentiating the slides. 100 mones are sufficient to separate frozen normal (green) from frozen tumor (orange) slides. (b) Venn diagram of statistically significant mones differentiating tumor from adjacent normal frozen slides, comparing different statistical tests. Venn diagrams were calculated for each cancer type, and the observed plot shows the average across all cancer types. On average the statistical tests agree on 75% of mones differentiating between tumor and normal slides. (c) Probability density function of mone 983 among frozen tumor (orange) and adjacent normal (green) BRCA slides. (d) Log-normalized scatter plot of slide level mone 983 and cellpose (see [53]) estimates of cellularity across BRCA frozen slides. (e) Example tiles from slides with extreme mone 983 values (high and low). (f) Cluster map of the mone-mone correlation matrix of LUAD tumor slides, demonstrating that many mones pairs are highly correlated. (g) Mone-mone correlation matrix of LUSC slides, with mones ordered identically to the Figure 1f cluster map.

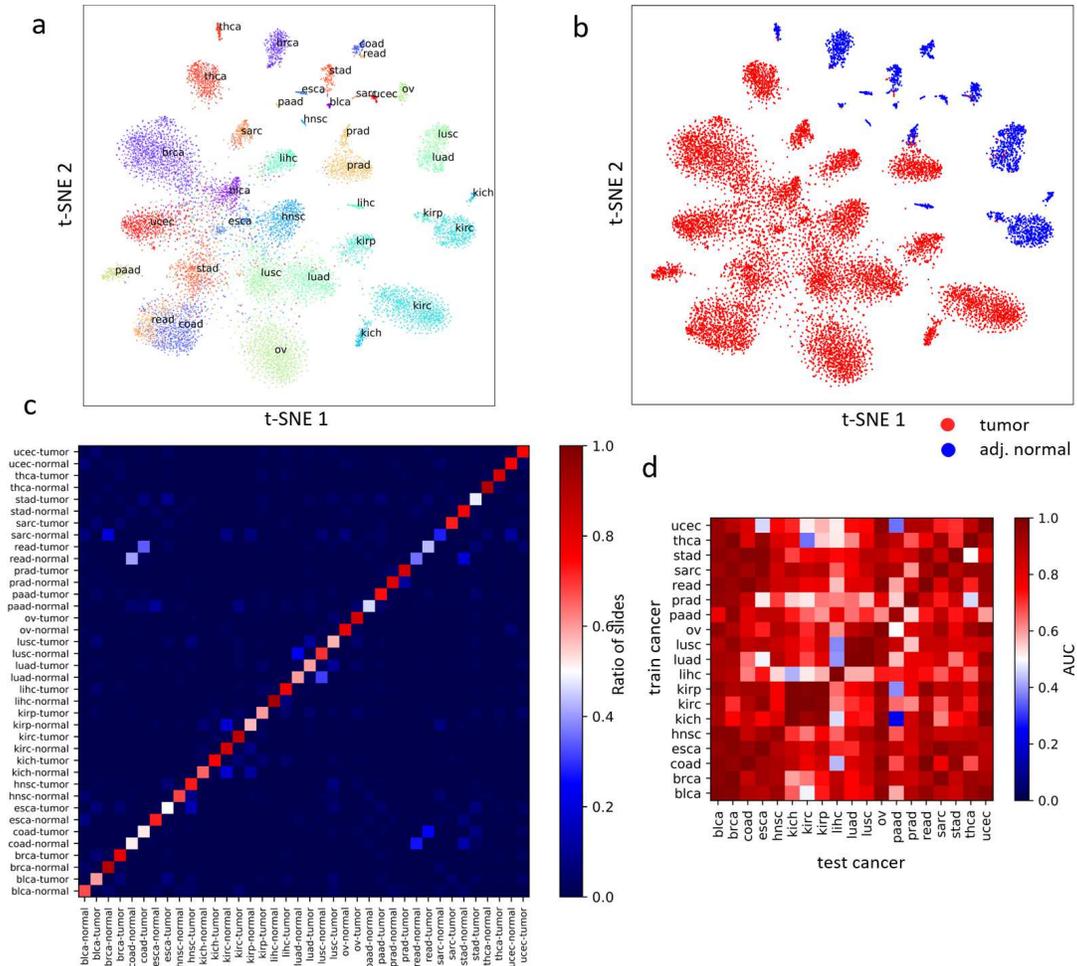


Figure 3: The joint distribution of mones reliably separates tumor and normal slides and the underlying cancer. 2D t-SNE plots of the mones-based MLDA feature space distinguishing 38 classes (19 cancers, tumor/normal status) based on (a) cancer type and (b) tumor/normal status. (c) Normalized confusion matrix of the 38-class mones-based logistic regression classifier. The color depicts the ratio of slides with a given true class predicted as any of the possible classes. The large diagonal values suggest the classifier has high accuracy. (d) The cross-classification AUCs of mones-based logistic regression tumor/normal classifier trained on each cancer and applied to all cancers.

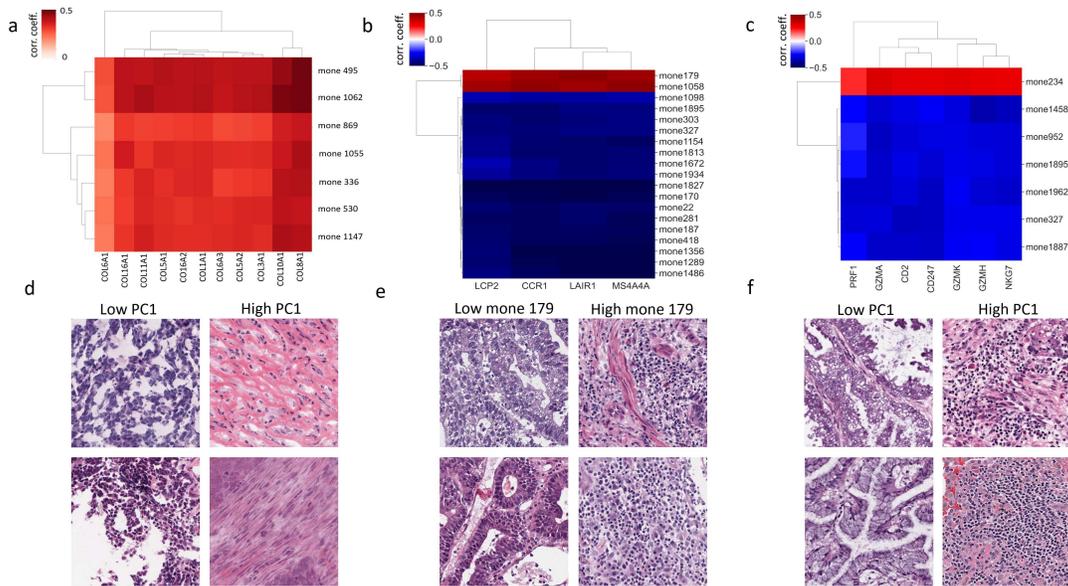


Figure 4: Mone-gene correlation analysis identifies highly correlated mone-gene clusters. Correlation matrix of (a) a cluster of highly correlated mones and collagen genes in OV, and a cluster of highly correlated mones and immune-related genes in (b) COAD and (c) LUAD. See Supplementary Figure 9 for adjusted p-values. Example tiles from slides with (d) high and low PC-1 in OV, (e) high and low mone 179 in COAD, and (f) high and low PC-1 in LUAD. Histopathology review identifies that mone-predicted (d) OV tiles with high PC-1 are rich in collagen, and (e) COAD with high mone 179 and (f) luad tiles with high PC-1 have a strong lymphocyte presence. See Supplementary Figures 10-13 for additional examples at both high and low mone values.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryFigures.pdf](#)
- [SupplementaryFiles.zip](#)